



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Bayesian geophysical inversion using invertible neural networks

Citation for published version:

Zhang, X & Curtis, A 2021, 'Bayesian geophysical inversion using invertible neural networks', *Journal of Geophysical Research. Solid Earth*, vol. 126, no. 7, e2021JB022320.
<https://doi.org/https://agupubs.onlinelibrary.wiley.com/doi/10.1029/2021JB022320?af=R>

Digital Object Identifier (DOI):

<https://agupubs.onlinelibrary.wiley.com/doi/10.1029/2021JB022320?af=R>

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Journal of Geophysical Research. Solid Earth

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Bayesian geophysical inversion using invertible neural networks

Xin Zhang¹ and Andrew Curtis¹

¹School of Geosciences, University of Edinburgh, Edinburgh, United Kingdom

Key Points:

- We introduce invertible neural networks that solve Bayesian geophysical inverse problems probabilistically.
- We use the networks to solve two types of problems: surface wave dispersion inversion and travel time tomography.
- For repeated inverse problems the method provides an efficient and accurate approximation to results obtained using Monte Carlo sampling.

Corresponding author: Xin Zhang, xzhang15@ed.ac.uk

Abstract

Constraining geophysical models with observed data usually involves solving nonlinear and non-unique inverse problems. Mixture density networks (MDNs) provide an efficient way to estimate Bayesian posterior probability density functions (pdf's) that represent the non-unique solution. However it is difficult to infer correlations between parameters using MDNs, and in turn to draw samples from the posterior pdf. We introduce an alternative to resolve these issues: invertible neural networks (INNs). These are simultaneously trained to represent uncertain forward functions and to solve Bayesian inverse problems. In its usual form, the method does not account for uncertainty caused by data noise and becomes less effective in high dimensionality. To overcome these issues, in this study we include data noise as additional model parameters and train the network by maximising the likelihood of the data used for training. We apply the method to two types of imaging problems: 1D surface wave dispersion inversion and 2D travel time tomography, and compare the results to those obtained using Monte Carlo and MDNs. Results show that INNs provide comparable posterior pdfs to those obtained using Monte Carlo, including correlations between parameters, and provide more accurate marginal distributions than MDNs. After training, INNs estimate posterior pdfs in seconds on a typical desktop computer. Hence they can be used to provide efficient solutions for repeated inverse problems using different data sets. Even accounting for training time, our results also show that INNs can be more efficient than Monte Carlo methods for solving inverse problems only once.

1 Introduction

Geoscientists build models of the subsurface in order to understand properties and processes in the Earth's interior. The models are usually parameterized in some way, so to constrain the models we must solve a parameter estimation problem. Data are recorded which provide constraints, together with prior knowledge. However, since the physical relationships between parameters and data usually predict data given the parameters (known as the forward calculation) but not the reverse, the solution must be found using inverse theory.

Geophysical inverse problems usually have non-unique solutions due to noise in the data, to nonlinearity of the physical relationships between model parameters and data, and to fundamentally unconstrained combinations of parameters. Uncertainties in pa-

parameter estimates must therefore be quantified in order to interpret inversion results correctly. Unfortunately, estimating uncertainty in nonlinear inverse problems can be computationally expensive, and the cost increases both with the number of parameters, and with the computational cost of the forward calculation. In this study we therefore solve two different types of seismic tomography problems which each have fewer than 100 parameters, and have relatively rapid forward functions (each evaluation takes on the order of seconds). This allows us to evaluate solutions sufficiently accurately to thoroughly test a new method of Geophysical inversion.

Geophysical inverse problems are traditionally solved by linearising (approximating) the nonlinear physics, and using optimization methods which seek a model that minimizes the misfit between observed and predicted data (Aki & Lee, 1976; Dziewonski & Woodhouse, 1987; Iyer & Hirahara, 1993). However, despite their wide applications, linearised procedures cannot produce accurate estimates of uncertainty because data uncertainty distributions can deviate substantially from analytic forms of probability such as Gaussians or Uniform distributions that can be propagated through linearised methods efficiently, and since the distribution of possible models post inversion can be strongly affected by nonlinearity in the physics (Bodin & Sambridge, 2009; Smith, 2013; Galetti et al., 2015; Zhang et al., 2018). Methods based on nonlinear Bayesian formulations of inverse problems have been introduced to provide more accurate uncertainty estimates. In such methods one finds a so-called probability density function (pdf) which describes the non-uniqueness of the parameters, by combining information independent of the data called *prior* information with information in the data. These methods include Markov chain Monte Carlo (MCMC) sampling methods (Mosegaard & Tarantola, 1995; Sambridge, 1999; Malinverno et al., 2000; Bodin & Sambridge, 2009; Galetti et al., 2015; Zhang et al., 2018) and variational inference (M. A. Nawaz & Curtis, 2018; M. Nawaz & Curtis, 2019; M. A. Nawaz et al., 2020; Zhang & Curtis, 2020a, 2020b).

MCMC methods generate a set of samples from the posterior probability density function (Brooks et al., 2011), which can be used thereafter to derive useful statistics which describe that pdf (e.g. mean, standard deviation, etc.). MCMC methods are quite general from a theoretical point of view and have been applied to a range of geophysical inverse problems, for example, to surface wave dispersion inversion (Bodin et al., 2012; Shen et al., 2012; Young et al., 2013; Galetti et al., 2017; Zhang et al., 2018; Zhang, Hansteen, et al., 2020), seismic travel time tomography (Bodin & Sambridge, 2009; Galetti et al.,

2015; Piana Agostinetti et al., 2015; Fichtner et al., 2018; Zhang, Roy, et al., 2020), full-waveform inversion (Ray et al., 2016, 2017; Gebraad et al., 2020; Khoshkholgh et al., 2020; Kotsi et al., 2020), gravity inversion (Bosch et al., 2006; Rossi, 2017), magnetotellurics inversion (Jones & Hutton, 1979), electromagnetic inversion (Minsley, 2011; Ray et al., 2013) and fluid flow history matching (Subbey et al., 2003; Mohamed et al., 2010). However, such solutions are acquired at significant expense, typically requiring days or weeks of computer run time, and hence cannot be applied in scenarios that require rapid solutions such as real-time monitoring (Duputel et al., 2009; Cao et al., 2020), or when many similar inversions must be performed (Käufel et al., 2016).

Variational inference provides a different way to solve Bayesian inference problems. The method seeks an optimal approximation to the posterior pdf within a predefined, expressive family of probability distributions by minimizing the Kullback-Leibler divergence between the approximating pdf and the posterior pdf (Blei et al., 2017). Since the method solves the inference problem using optimization rather than stochastic sampling, it can be more computationally efficient than Monte Carlo methods. Variational methods have been applied to invert for geological facies and petrophysical properties (M. A. Nawaz & Curtis, 2018; M. Nawaz & Curtis, 2019; M. A. Nawaz et al., 2020), and for travel time tomography (Zhang & Curtis, 2020a; Zhao et al., 2020), full waveform inversion (Zhang & Curtis, 2020b) and seismic denoising (Siahkoobi et al., 2020). However, although variational inference can be relatively efficient, it still typically requires large compute times to obtain solutions for the problems, and therefore may not provide sufficiently rapid solutions for real-time monitoring, nor for cases where many similar inversions are required.

Neural network based methods offer another efficient alternative for certain classes of inverse problems that must be solved many times with new data of similar type. An initial set of Monte Carlo samples is generated from the pdf that describes the *a priori* information (the so-called prior pdf), and synthetic data are simulated computationally for each of these samples. Neural networks are flexible mappings which can be trained to emulate any specific inverse mapping from data to parameter space by fitting this set of examples of that mapping (called the training data set; Bishop, 2006). Thereafter the trained neural networks interpolate the inverse mapping between the training examples, and can be evaluated efficiently for any new, measured data set to provide estimates of corresponding parameter values. Hence they can be applied to applications that require solutions to many different inverse problems within the class of problems represented by

the training data. Neural networks were first introduced to Geophysics by Röth and Taramola (1994) to estimate subsurface velocity from active source seismic waveforms, and have been applied to seismic velocity inversion using earthquake data (Moya & Irikura, 2010) and semblance gathers (Araya-Polo et al., 2018) and to seismic anisotropy inversion (You et al., 2020). Laloy et al. (2019) introduced vector-to-image transfer networks to solve inverse problems and applied them to transient groundwater flow and ground penetrating radar tomographic problems. Mosser et al. (2020) used so-called generative adversarial networks to re-parameterize geologically correlated Earth structure with a relatively low number of parameters, and inverted for the structure that best fit synthetic seismic waveform data.

The above studies did not provide estimates of uncertainties since for each input data vector their neural networks only predict one parameter vector. Devilee et al. (1999) proposed the first geophysical probabilistic form of neural networks which provide discretised Bayesian posterior pdfs, and used them to invert surface wave dispersion data for crustal thickness maps and their uncertainties across Eurasia. In an alternative formulation, mixture density networks (MDNs) output a probability distribution that is defined by a sum of analytic pdfs called kernels, such as Gaussian distributions, and can be trained to map data to corresponding posterior pdfs (Bishop, 2006). MDNs have been applied to surface wave dispersion inversion (Meier et al., 2007b, 2007a; Earp et al., 2020; Cao et al., 2020), 2D travel time tomography (Earp & Curtis, 2020), petrophysical inversion (Shahraeeni & Curtis, 2011; Shahraeeni et al., 2012), earthquake source parameter estimation (Käufel et al., 2014, 2015), Earth’s radial seismic structure inversion (de Wit et al., 2013), pore pressure prediction (Karmakar & Maiti, 2019), mapping of lithology (Karmakar et al., 2018), wind prediction (Men et al., 2016), acoustic-articulatory inversion (Richmond, 2007) and nuclei detection (Koohababni et al., 2018). However MDNs become difficult to train in high dimensionality because of numerical instability, and they suffer from mode collapse which means that some modes of the posterior pdf are missing in the results (Hjorth & Nabney, 1999; Rupprecht et al., 2017; Curro & Raquet, 2018; Cui et al., 2019; Makansi et al., 2019). Consequently they are less effective at inferring correlations between parameters, so in practice usually very low (often single) dimensional marginal distributions are inferred (Meier et al., 2007b, 2007a; Earp & Curtis, 2020; Earp et al., 2020).

To estimate full posterior pdfs, Ardizzone et al. (2018) proposed to use invertible neural networks (INNs) to solve probabilistic inverse problems, and showed that the method can provide accurate approximations to multimodal, highly correlated posterior distributions. INNs provide bijective mappings between inputs (models) and outputs (data), and can be trained to estimate posterior pdfs by introducing additional latent variables in the outputs (data) side. They have been used to solve inverse problems in medicine (Ardizzone et al., 2018), astrophysics (Osborne et al., 2019), optical imaging (Adler et al., 2019; Moran et al., 2018) and morphology (Sahin & Gurevych, 2020). In this study we use INNs to solve seismic tomographic inverse problems. Note that INNs have also been used to solve a variational problem to parameterize uncertainty for reservoir characterization (Rizzuti et al., 2020). The idea of auxiliary variables has also been used in seismic full waveform inversion (Huang et al., 2018).

The INN method proposed in Ardizzone et al. (2018) does not account for uncertainty caused by data noise in the results. To resolve this issue, in this study we include the random data noise as additional model parameters. In addition, the method uses maximum mean discrepancy to measure differences between two distributions during training, and varies the network parameters so as to minimise this measure. However, this measure becomes less effective as the dimensionality increases because of the curse of dimensionality (Ramdas et al., 2015). We show that this issue can be overcome by using a maximum likelihood criterion to train INNs.

In the next section we describe the basic structure of INNs, and how they can be trained to solve Bayesian inference problems. We then apply the method to two types of seismic inverse problems: 1D surface wave dispersion inversion and 2D travel time tomography, and compare the results with those obtained using Markov chain Monte Carlo (McMC) and MDNs. We demonstrate that INNs can provide comparable probabilistic results with those obtained using McMC, including correlations between parameters, whereas MDNs provide far less information about inter-parameter correlations. In our travel time tomography examples the computational time of training INNs and MDNs, including generation of the synthetic training data is comparable to one single run of McMC. We thus demonstrate that INNs can provide fast, accurate approximations of posterior pdfs even if the problem is solved only once; they can then produce rapid solutions for subsequent problems within the same problem class.

2 Methods

2.1 Bayesian inference

Bayesian methods update a *prior* probability density function (pdf) $p(\mathbf{m})$ with new information from data \mathbf{d}_{obs} to produce a probability distribution of model parameters \mathbf{m} post inversion, which is often called a *posterior* pdf and written as $p(\mathbf{m}|\mathbf{d}_{obs})$. According to Bayes' theorem,

$$p(\mathbf{m}|\mathbf{d}_{obs}) = \frac{p(\mathbf{d}_{obs}|\mathbf{m})p(\mathbf{m})}{p(\mathbf{d}_{obs})} \quad (1)$$

where $p(\mathbf{d}_{obs}|\mathbf{m})$ is the *likelihood* which is the probability of observing data \mathbf{d}_{obs} if model \mathbf{m} was true, and $p(\mathbf{d}_{obs})$ is a normalization factor called the *evidence*. The likelihood function is often assumed to follow a Gaussian probability density function around the data predicted synthetically (using known, usually nonlinear physical relationships) from model \mathbf{m} , as this is assumed to be a reasonable approximation to the pdf of uncertainties or errors in measured data. Estimating the posterior distribution given prior information and the likelihood is called Bayesian inference.

2.2 Invertible neural networks

Invertible neural networks (INN) are a class of networks that provide bijective mappings between inputs and outputs. A typical design of an INN contains a serial sequence of reversible blocks, each of which consists of two coupled layers (Dinh et al., 2016; Kingma & Dhariwal, 2018). Each block's input vector \mathbf{u} (for example, \mathbf{u} can be the model vector) is split into two halves \mathbf{u}_1 and \mathbf{u}_2 , which are transformed by an affine function with coefficients $exp(s_i)$ and t_i to produce the output $[\mathbf{v}_1, \mathbf{v}_2]$:

$$\begin{aligned} \mathbf{v}_1 &= \mathbf{u}_1 \odot exp(s_2(\mathbf{u}_2)) + t_2(\mathbf{u}_2) \\ \mathbf{v}_2 &= \mathbf{u}_2 \odot exp(s_1(\mathbf{u}_1)) + t_1(\mathbf{u}_1) \end{aligned} \quad (2)$$

where \odot represents element-wise multiplication. This process is trivially invertible for any functions t and s :

$$\begin{aligned} \mathbf{u}_2 &= (\mathbf{v}_2 - t_1(\mathbf{u}_1)) \odot exp(-s_1(\mathbf{u}_1)) \\ \mathbf{u}_1 &= (\mathbf{v}_1 - t_2(\mathbf{u}_2)) \odot exp(-s_2(\mathbf{u}_2)) \end{aligned} \quad (3)$$

Importantly functions s_i and t_i do not need to be invertible themselves. In this study we use fully connected neural networks or convolutional neural networks to represent trainable functions s_i and t_i . To improve interaction between parameters of the input vec-

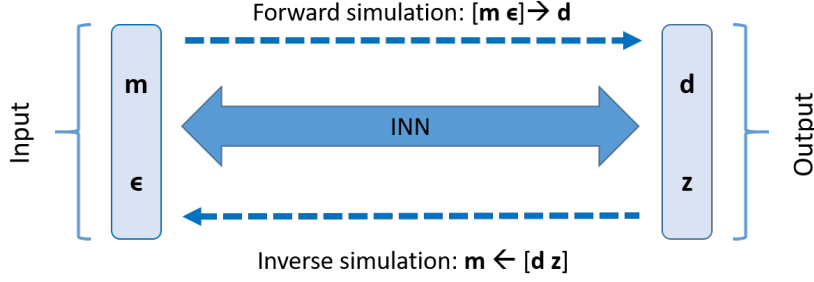


Figure 1. A conceptual figure of invertible neural networks. A latent random variable \mathbf{z} is added to the outputs to account for uncertainties in the inputs \mathbf{m} . The latent variable can follow any probability distribution, and is chosen to follow a standard Gaussian distribution in this study. The posterior distribution of \mathbf{m} can be obtained by sampling \mathbf{z} for a fixed measurement \mathbf{d} and running the trained neural network backwards. To appropriately account for noise in the data, we include random data noise ϵ as additional model parameters.

tor, we add a permutation layer after each reversible block (equation 2), which shuffles outputs of that block in a randomized, but fixed way as in Ardizzone et al. (2018).

2.3 Solving inverse problems using INNs

We first review the idea of using INNs to solve inverse problems as proposed by Ardizzone et al. (2018). INNs provide a natural way to solve inverse problems. For example, training an INN in the forward direction on a well-understood forward process $\mathbf{d} = F(\mathbf{m})$, one can obtain a solution to the inverse problem for free by running the trained network in the reverse direction. However, in practice inverse problems often have nonunique solutions. To account for uncertainties in the solution, additional latent variables \mathbf{z} (Figure 1) can be introduced to the outputs \mathbf{d} (Ardizzone et al., 2018). Though multiple models \mathbf{m} can produce the same \mathbf{d} , the pair comprising one model \mathbf{m} and the augmented vector $[\mathbf{d}, \mathbf{z}]$ is unique. So the purpose of adding latent variables \mathbf{z} is that the relationship between the range of models and the space of augmented vectors becomes one-to-one. The networks therefore associate model parameters \mathbf{m} to a unique pair $[\mathbf{d}, \mathbf{z}]$ of measurements and latent variables, written as $[\mathbf{d}, \mathbf{z}] = f(\mathbf{m}; \theta)$ where θ represents parameters of the neural networks. As in Ardizzone et al. (2018) one can train the neural network to approximate the forward process, that is $f(\mathbf{m}; \theta)_d \approx F(\mathbf{m})$ where the subscript d represents the data part of the network output, and meanwhile ensure the latent vari-

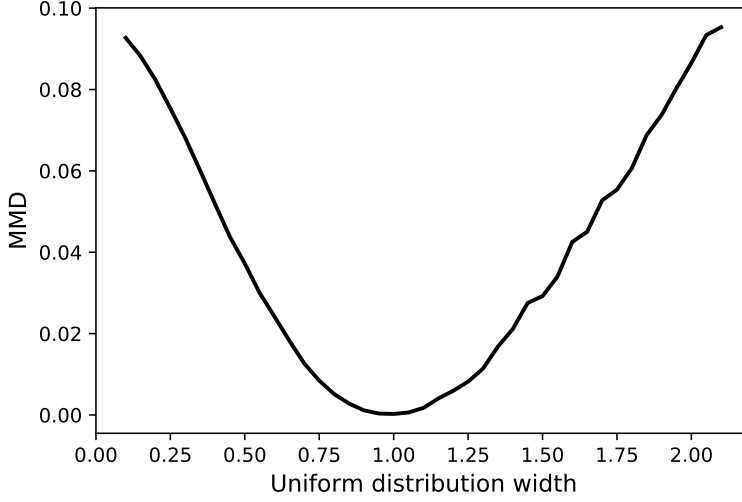


Figure 2. An example of MMD between two Uniform distributions which have the same mean value but different widths: one has a fixed width 1, the other one has various widths from 0.1 to 2.1 (horizontal axis in the figure). MMD reaches zero when the two distributions are the same.

able \mathbf{z} predicted by the network are distributed according to a chosen distribution, for example, a Gaussian distribution. The solution of the inverse problem can be obtained thereafter by running the network backwards given a specific measurement \mathbf{d}^* with latent variable \mathbf{z} selected randomly from the same Gaussian distribution:

$$\begin{aligned} \mathbf{m} &= f^{-1}(\mathbf{d}^*, \mathbf{z}; \theta) \\ \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{aligned} \tag{4}$$

By taking many samples of $[\mathbf{d}^*, \mathbf{z}]$ where \mathbf{z} is generated from the chosen Gaussian distribution and \mathbf{d}^* is repeated, the trained network transforms the distribution $p(\mathbf{z})$ to the posterior distribution $p(\mathbf{m}|\mathbf{d}^*)$, which describes the distribution of models that can generate the exact data \mathbf{d}^* (we introduce data measurement uncertainties below). Though in equation 4 we used a Gaussian distribution for the latent variable \mathbf{z} , theoretically any other distributions can be used. Define the output distribution of the network in the forward direction to be $q(\mathbf{d}, \mathbf{z}; \theta)$:

$$q(\mathbf{d}, \mathbf{z}; \theta) = p(\mathbf{m})/|\det J_f(\mathbf{m}; \theta)| \tag{5}$$

where $p(\mathbf{m})$ is the prior distribution of model \mathbf{m} , $J_f(\mathbf{m}; \theta) = \frac{\partial f(\mathbf{m}; \theta)}{\partial \mathbf{m}}$ is the Jacobian of the forward transform embodied in the network (Devore & Berk, 2012). Given those

expressions, the training loss function \mathcal{L} can be expressed as:

$$\mathcal{L} = \sum_{i=1}^N \|\mathbf{d}^i - f(\mathbf{m}^i; \theta)\| + \alpha \text{MMD}[q(\mathbf{d}^i, \mathbf{z}^i; \theta), p(\mathbf{d}^i)p(\mathbf{z}^i)] \quad (6)$$

where the superscript i denotes the index of batches and N is the number of batches. $p(\mathbf{d}^i)$ is the prior distribution of data of the i^{th} batch which is generated by applying the forward function $F(\mathbf{m})$ over the prior distribution $p(\mathbf{m}^i)$, MMD represents the Maximum Mean Discrepancy which is a measure of difference between two distributions, and α is the relative weight of the MMD term. Note that this loss function only needs to be trained in the forward direction. MMD can be evaluated using only samples from the two distributions in its arguments. For example, assume X and X' are random variables with distribution p , and that Y and Y' are random variables with distribution q , then the MMD can be expressed as:

$$\text{MMD}[p, q] = \text{E}_{X, X'}[k(X, X')] - 2\text{E}_{X, Y}[k(X, Y)] + \text{E}_{Y, Y'}[k(Y, Y')] \quad (7)$$

where k is a kernel function. Here we use the Inverse Multiquadratic function $k(x, x') = 1/(1 + \|(x - x')/h\|_2^2)$ as it has heavier tails than a Gaussian kernel and provides meaningful gradients for outliers (Tolstikhin et al., 2017). The MMD equals zero if and only if $p = q$ (Gretton et al., 2012). Figure 2 shows an example of MMD between two Uniform distributions which have the same mean value but different widths: one has a fixed width 1, while the other has various widths from 0.1 to 2.1. The MMD reaches zero when the two distributions are the same. Ardizzone et al. (2018) showed that when the loss in equation (6) reaches zero, the neural network produces the posterior distribution $p(\mathbf{m}|\mathbf{d}^*)$ for a given measurement \mathbf{d}^* (see proof of this result in Appendix A). In practice to facilitate the convergence of training, a loss function is also included on the input side:

$$\mathcal{L} = \sum_{i=1}^N \|\mathbf{d}^i - f(\mathbf{m}^i; \theta)\| + \alpha \text{MMD}[q(\mathbf{d}^i, \mathbf{z}^i; \theta), p(\mathbf{d}^i)p(\mathbf{z}^i)] + \beta \text{MMD}[q(\mathbf{m}^i; \theta), p(\mathbf{m}^i)] \quad (8)$$

where $q(\mathbf{m}^i; \theta) = p(\mathbf{d}^i)p(\mathbf{z}^i)/|J_{f^{-1}}(\mathbf{d}^i, \mathbf{z}^i; \theta)|$ is the input distribution predicted by the neural network acting in the inverse direction and β is the relative weight. In the training process the hyperparameters α and β are selected manually to give a balance between different terms and are kept the same for all batches. To train the network, the parameters θ are optimized using the ADAM method of stochastic gradient descent (Kingma & Ba, 2014) to minimize the loss function in equation 8. After training, the network provides an approximation function f which maps a model to its corresponding data, as well

as a distribution $q(\mathbf{d}, \mathbf{z})$ which provides an estimate of the true distribution $p(\mathbf{d})p(\mathbf{z})$. For a specific measurement \mathbf{d}^* , the posterior distribution is obtained using equation 4.

2.4 Bayesian inference using INNs

The method described above only accounts for intrinsic uncertainties caused by non-linearity of underlying physics and neglects uncertainties caused by data noise. This is because theoretically the method requires the likelihood function to be a delta function which cannot be achieved in practice due to noise in the data (see the proof in Appendix A). To appropriately account for noise in the data such that the estimated posterior pdfs using INNs are consistent with the posterior pdfs in Bayesian inference, we treat random noise as additional model parameters (Figure 1), that is, we assume:

$$\begin{aligned}\mathbf{d} &= F(\mathbf{m}) + \epsilon \\ \epsilon &\sim \mathcal{N}(\mathbf{0}, \Sigma)\end{aligned}\tag{9}$$

where Σ is the covariance matrix of data noise. Although here we assume that noise follows a Gaussian distribution, in principle any other noise distributions can be used. Note that this idea of treating noise as additional parameters has also been used in training MDNs (Earp et al., 2020), but in that case the inputs to the network were the standard deviations of measured data uncertainties. Although in principle we could include a full covariance matrix, in this study we assume Σ to be a diagonal matrix. The input dimensionality of the network now becomes (see Figure 1):

$$\dim(\text{inputs}) = \dim(\mathbf{m}) + \dim(\mathbf{d})\tag{10}$$

In this way the estimated posterior pdfs approximate the correct solution for Bayesian inference (see discussion in Appendix A). After training for a specific, noisy measurement \mathbf{d}_{obs} , the Bayesian posterior pdf in equation 1 can be obtained similarly using equation 4 with \mathbf{d}^* substituted with \mathbf{d}_{obs} . Since INNs require the same dimensionality of inputs and outputs, the dimensionality of \mathbf{z} becomes:

$$\begin{aligned}\dim(\mathbf{z}) &= \dim(\text{outputs}) - \dim(\mathbf{d}) \\ &= \dim(\text{inputs}) - \dim(\mathbf{d}) \\ &= \dim(\mathbf{m})\end{aligned}\tag{11}$$

And because the network also needs to capture the distribution of noise parameter ϵ , the dimensionality of \mathbf{z} will need to be higher than $\dim(\mathbf{m})$. In this case zeros can be padded to the input side to match the input and output dimensionalities.

Note that trained INNs also provide approximating forward functions. For example one can obtain the distribution of data with noise by running the network forward with noise parameter ϵ distributed according to its assumed distribution given a fixed model \mathbf{m} . In our case since we assumed Gaussian noise, the same distribution of data can actually be obtained by adding random noise to the synthetic data. However in cases in which noise distributions are not explicitly known, for example noise caused by assumptions in forward modelling, INNs provide a way to generate associated data distributions.

Although MMD is an efficient method to discriminate between two distributions in low dimensionality, it becomes less efficient (requires many more samples) in high dimensionality (Ramdas et al., 2015). To improve the efficiency of the method in such cases we add a maximum likelihood term to the loss function. Define the probability distribution of $[\mathbf{m}, \epsilon]$ as $p(\mathbf{m}, \epsilon) = p(\mathbf{m})p(\epsilon)$ where $p(\epsilon)$ is the distribution of noise parameter ϵ . The distribution of the network output $[\mathbf{d}, \mathbf{z}]$ can be expressed as:

$$\begin{aligned} q(\mathbf{d}, \mathbf{z}; \theta) &= p(\mathbf{m}, \epsilon) |det J_{f^{-1}}(\mathbf{d}, \mathbf{z}; \theta)| \\ &= p\left[\left(\mathbf{m}, \epsilon\right) = f^{-1}(\mathbf{d}, \mathbf{z}; \theta)\right] |det J_{f^{-1}}(\mathbf{d}, \mathbf{z}; \theta)| \end{aligned} \quad (12)$$

where we have used the fact that $[\mathbf{m}, \epsilon] = f^{-1}(\mathbf{d}, \mathbf{z})$. The distribution in equation 12 is the likelihood function of data $[\mathbf{d}, \mathbf{z}]$ given the network parameters θ . Thus, we can train the network by maximizing the likelihood function:

$$\begin{aligned} \mathcal{L} &= \sum_{i=1}^N \|\mathbf{d}^i - f(\mathbf{m}^i, \epsilon^i; \theta)\| + \alpha \text{MMD}[q(\mathbf{d}^i, \mathbf{z}^i; \theta), p(\mathbf{d}^i)p(\mathbf{z}^i)] + \beta \text{MMD}[q(\mathbf{m}^i, \epsilon^i; \theta), p(\mathbf{m}^i)p(\epsilon^i)] \\ &\quad - \gamma \log(p\left[\left(\mathbf{m}^i, \epsilon^i\right) = f^{-1}(\mathbf{d}^i, \mathbf{z}^i; \theta)\right] |det J_{f^{-1}}(\mathbf{d}^i, \mathbf{z}^i; \theta)|) \end{aligned} \quad (13)$$

where γ is the relative weight of the likelihood term. Maximising this likelihood term with respect to the parameters θ of the network for all training data (\mathbf{d}, \mathbf{z}) ensures that the network transforms between distributions $p(\mathbf{m})p(\epsilon)$ and $p(\mathbf{d})p(\mathbf{z})$ (Bishop, 2006). Note that this likelihood term therefore achieves the same goal as the MMD terms alone, but is more effective in high dimensionality. We demonstrate the efficiency of the new loss function using our second example in the online supporting information.

Figure 3 shows a toy example application of the method. The training data (Figure 3a) are generated using a function $y = x^2 + \epsilon$ where $x \sim \text{Uniform}(-1, 1)$ and $\epsilon \sim \mathcal{N}(0, 0.04)$. We train INNs to predict the posterior pdf $p(x|y^*)$ for a given y^* using methods described in section 2.3 and 2.4. For example, Figure 3b shows the pdf predicted by the trained INN using the method in section 2.4 (with ϵ) at $y = 0.6$ (orange histogram), which provides an accurate approximation to the results obtained by Markov chain Monte Carlo (blue line). In comparison Figure 3c shows the pdf at $y = 0.6$ predicted by the trained INN using the method in section 2.3 (without ϵ). Figure 3d shows the zoomed-in distributions obtained using different methods. The results show that without adding noise parameters to the modelling process, the trained INN can produce biased results compared to McMC. Figure 3e shows the distribution of ϵ (orange histogram) predicted by the trained INN (with ϵ), which matches the true distribution (blue line) as expected. In Figure 3d we show the data distribution predicted by the trained INN (with ϵ) when $x = 0.5$, which gives an accurate approximation to the true distribution (blue line). This example shows the potential of INNs to predict accurate posterior pdfs as well as to approximate probabilistic forward functions.

3 Synthetic tests

3.1 1D surface wave dispersion inversion

As a first experiment we train an INN to predict posterior pdfs for 1D seismic velocity structure with depth, given surface wave dispersion data. The subsurface is parameterized using ten regular layers with 0.1 km spacings for the shallower part ($< 0.5\text{km}$) and 0.2 km spacings for the deeper part ($> 0.5\text{km}$) since surface wave dispersion inversions are known to suffer diminishing spatial resolution with depth. For each layer we specify a Uniform prior distribution for the shear wave velocity (Figure 4a), whose velocity range is set to be typical for the near surface (Zhang, Hansteen, et al., 2020). P-wave velocity and density are calculated from the shear velocity using $V_p = 1.16V_s + 1.36$ (Castagna et al., 1985) and $\rho = 1.74V_p^{0.25}$ (Brocher, 2005) where V_p and V_s are P and S wave velocity in km/s and ρ is density in g/cm^3 . We generate 100,000 models from the prior pdf and calculate Rayleigh wave dispersion curves corresponding to each model using a modal approximation method (Herrmann, 2013) over the period range 0.7 s to 2.0 s with 0.1 s spacing (Figure 4b). We added Gaussian noise with a standard deviation of 5 m/s to those calculated dispersion curves, which is a typical noise level in near

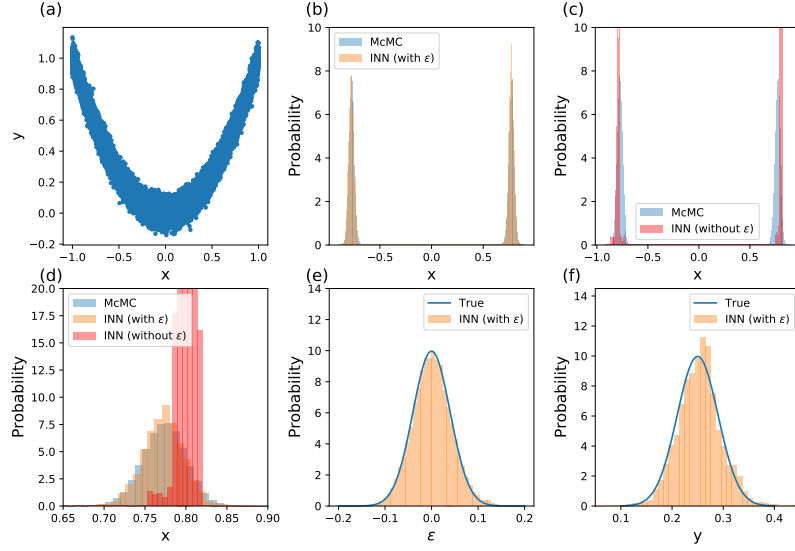


Figure 3. A toy example that uses INNs to predict posterior pdfs. **(a)** Training dataset. **(b)** Posterior pdfs of x obtained using INN with ϵ (orange histogram) compared with the results obtained using McMC (blue histogram) when $y = 0.6$. **(c)** Same posterior pdfs as in (b) but obtained using INN without ϵ (red histogram). **(d)** Zoomed-in posterior pdfs obtained using different methods. **(e)** Noise parameter distribution obtained using INN with ϵ (orange histogram) and the true distribution (blue line). **(f)** The distribution of y predicted by INN with ϵ (orange histogram) and the true distribution (blue line) when $x = 0.5$.

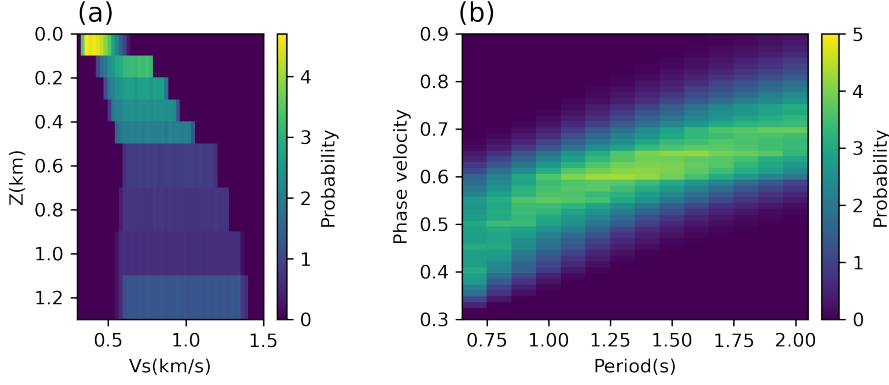


Figure 4. (a) Marginal prior distributions of shear velocities, and (b) the prior distribution of dispersion curves used to train neural networks.

surface ambient noise studies (Zhang, Hansteen, et al., 2020). Note that to ensure the computed dispersion curves are fundamental mode Rayleigh waves, within the prior pdf we ensured that the top layer has smallest shear velocity – otherwise the wave recorded on the Earth’s surface would be a higher mode Rayleigh wave (Zhang, Hansteen, et al., 2020). We use 90 percent of those model and dispersion curve pairs as training data, and the remaining 10 percent as test data used for independent evaluation of network performance.

The INN was designed using four reversible blocks, each of which contains fully connected subnetworks (see details in Appendix B1), and was trained using the ADAM optimizer (Kingma & Ba, 2014). We assume the neural network has converged when both the training loss and test loss become stationary or when the test loss starts increasing. The trained neural network is then used to predict posterior pdfs by running the network backwards using 5000 random values of \mathbf{z} in equation 4 with fixed data vector \mathbf{d}_{obs} , and we histogram the resulting set of samples of \mathbf{m} to approximate the posterior marginal distribution over each shear velocity. We show two examples of the set of predicted marginal posterior pdfs in Figure 5a and 5d. To better understand the results, we compared them with those obtained using McMC and MDNs. For McMC we use an adaptive Metropolis-Hastings algorithm (Haario et al., 2001; Salvatier et al., 2016) with 3 chains, each of which contains 250,000 samples including a burn-in period of 50,000 samples; the burn-in samples were ignored, and every 20th of the remaining 200,000 samples was included in the final set of McMC samples used for calculating statistics and marginal distributions. The

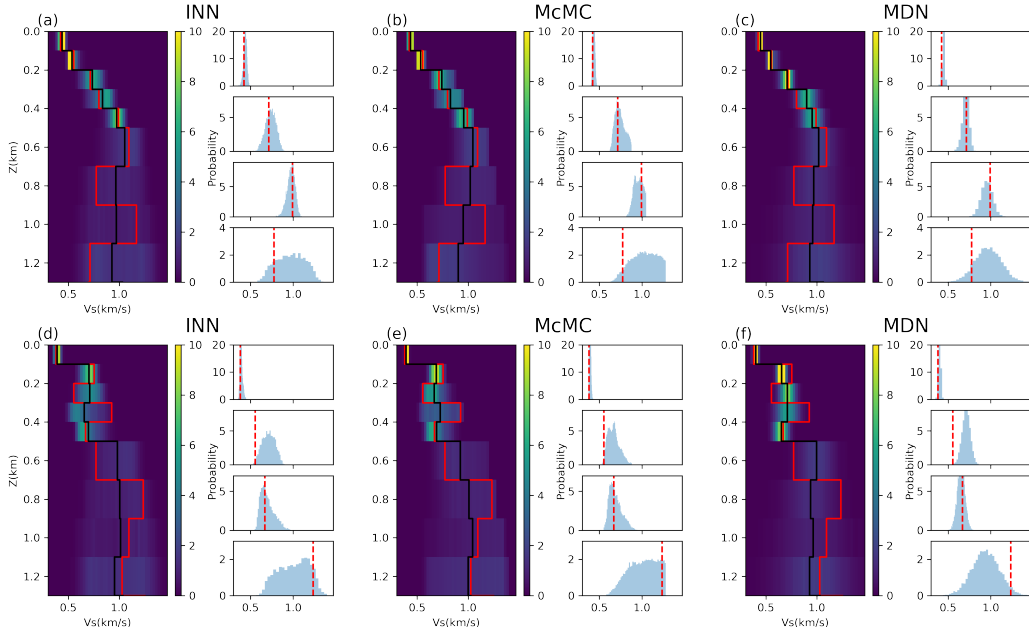


Figure 5. The marginal posterior distributions obtained using (a and d) an INN, (b and e) McMC and (c and f) MDNs for two different shear velocity structures. Red lines show the true shear velocity and black lines show the posterior mean velocity. At the right side of each panel we plot marginal distributions for four layers: 0 - 0.1 km, 0.2 - 0.3 km, 0.5 - 0.7 km and 0.9 - 1.1 km depth. Red dashed lines show the true velocity.

results are shown in Figure 5b and 5e. For the MDN we use 20 mixture Gaussian kernels and use a network design from Earp and Curtis (2020). The MDN network is then trained 20 times with random initialization and the network with best performance on the test data is used to produce final results. After training we generate 5000 samples from the MDN estimated posterior distribution and use them to calculate useful statistics and marginal distributions. The results are shown in Figure 5c and 5f.

Overall the three methods produce similar results for both examples. For example, in the top row the results show lower uncertainties at shallower depths ($< 0.8km$) due to the stronger prior information in the shallower part and the fact that surface waves are more sensitive to shallower structure. In the bottom row the results of all three methods exhibit higher uncertainties in the third and forth layers and the mean velocities deviate from the true velocity. This is because surface waves are not sensitive to complex structures with thin low velocity zones (Jan van Heijst et al., 1994). However, marginal distributions from the MDNs show clear Gaussian shapes, whereas the results from INNs and McMC have non-Gaussian shapes. This suggests that MDNs are not accurately approximating non-Gaussian pdfs, and in comparison the INNs have produced more accurate results. Note that this limitation of MDNs cannot be resolved by increasing the number of kernels as in these results only a few of the available kernels contribute to the final pdfs (all others are assigned near-zero weights), a property found also in previous studies (Hjorth & Nabney, 1999; Rupprecht et al., 2017; Curro & Raquet, 2018; Cui et al., 2019; Makansi et al., 2019; Earp & Curtis, 2020).

Figure 6 shows the correlation coefficients between parameters estimated using the three methods. The results from INNs and McMC show clear correlations between different parameters (Figure 6a,d and 6b,e), whereas the results from MDNs only show correlations for velocities of shallow layers in the top row. While this result probably occurs because we used a standard MDN which only contains kernels with diagonal covariance matrix, given the number of available kernels (20) it is certainly theoretically possible that the MDN could have represented the true correlations, at least approximately (in Figure 6c, the signs of correlations between the first, second and third layers are correct). Again, this relatively poor result is explained by the fact that MDNs tend to use only a few kernels to represent the solution (Hjorth & Nabney, 1999; Rupprecht et al., 2017; Curro & Raquet, 2018; Cui et al., 2019; Makansi et al., 2019; Earp & Curtis, 2020). Although a MDN with full covariance matrix might produce better results, the problem

becomes far more complex which can cause numerical instability and computationally expense due to the approximately squared number of network outputs that would be required (Williams, 1996).

In comparison INNs naturally estimate full correlation information. For example in our tests INNs provide the right correlation information for the first and second off-diagonal elements in the correlation matrix, that is, the correlation information between neighbouring and every second neighbouring layers (although the magnitudes are slightly lower than those from McMC). Even for more distant correlations, INNs still provide a reasonable approximation. For example, the results from INNs show correlations between the first and fourth layer that are similar to results from McMC. To further study correlation information between different parameters, in Figure 7 we show the bivariate marginal distributions of velocities at the second and fourth layer obtained using different methods for the two velocity profiles. The results from INN and McMC clearly show nonlinear trade-offs between the two parameters, whereas those obtained using MDN failed to find this information.

In the online supporting information we show the histograms of residuals (Figure S1) obtained using different methods for the two examples in Figure 5. The residuals are calculated using the ensemble of posterior samples obtained using different methods and are normalized by the noise level, so that the result should follow a standard Gaussian distribution $N(0, 1)$. The results show that the residual histograms obtained using INNs and McMC are close to the standard Gaussian distribution, whereas those obtained using MDN are different from the standard Gaussian distribution. This further demonstrates that INNs can provide more accurate posterior estimates than MDNs.

Overall INNs provide more accurate approximations to the results obtained using McMC compared to those obtained using MDNs. After training, both INNs and MDNs provide very efficient calculations of posterior distributions. For example, the above tests took about 2 seconds to predict posterior distributions using a trained INN or MDN on a typical desktop, whereas McMC took about 3 hours on the same machine. Thus trained INNs can be applied in scenarios where many repeated inversions are necessary to provide accurate shear velocity posterior distributions of the subsurface variations over space (Devilee et al., 1999; Meier et al., 2007a; Zhang, Hansteen, et al., 2020) or over time (Cao

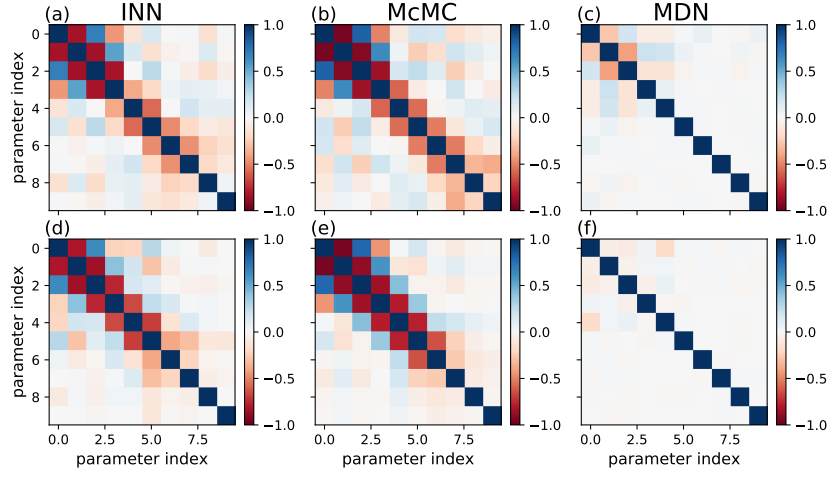


Figure 6. Correlation coefficients between shear velocities of different layers obtained using (a and d) an INN, (b and e) McMC and (c and f) MDNs for the two different velocity profiles in Figure 5.

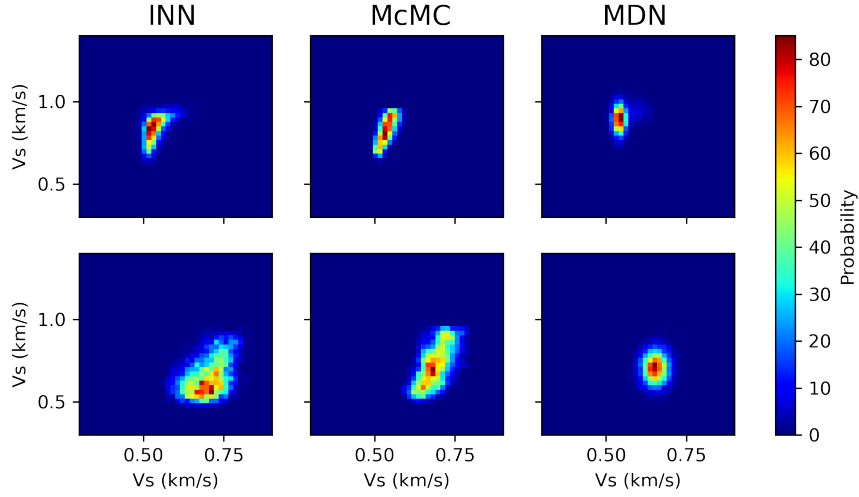


Figure 7. The bivariate marginal distributions of velocities at the second (horizontal axis) and fourth layer (vertical axis) obtained using INN, McMC and MDN respectively for the two velocity profiles in Figure 5.

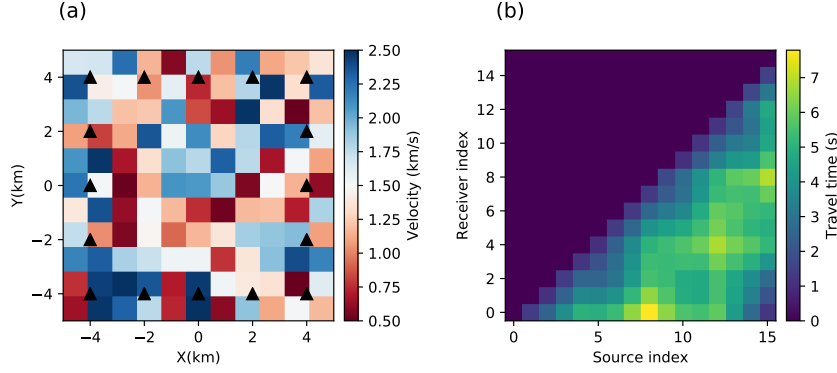


Figure 8. Experimental design for 2D travel time tomography. **(a)** Receiver locations (black triangles) and an example of a random velocity model. Each receiver also acts as a virtual source to mimic an ambient noise tomography experiment. **(b)** An example travel time field.

et al., 2020). We discuss the compute power required for training the INN in the Discussion section below.

3.2 Travel time tomography

We next examine performance of the methods on a 2D travel time tomographic problem similar to those which appear in ambient noise tomography (Shapiro et al., 2005). We solve a problem similar to that described in Earp and Curtis (2020) so that the results can be compared with theirs obtained using MDNs. A total of 16 receivers are used in our study, each of which also acts as a virtual source (Figure 8a). The velocity model is parametrized using a 9×9 regular grid (Figure 8a); for ease of visual interpretation we interpolate between the cell centres to construct smooth tomographic maps (e.g., Figure 9). At each grid point the prior distribution of velocity is set to be a Uniform distribution in the range 0.5 km/s to 2.5 km/s. We generate 200,000 velocity models from the prior distribution, for each of which the inter-receiver travel time data (Figure 8b) are calculated using a fast marching method (Rawlinson & Sambridge, 2004) with 0.05 s standard deviation Gaussian noise added. Note that we added a halo of cells with random velocities around the receiver array; these will not be imaged but do allow waves to travel both outside and inside of the array during inter-receiver propagation. Again we use 90 percent of those data as training data and the remaining 10 percent as test data.

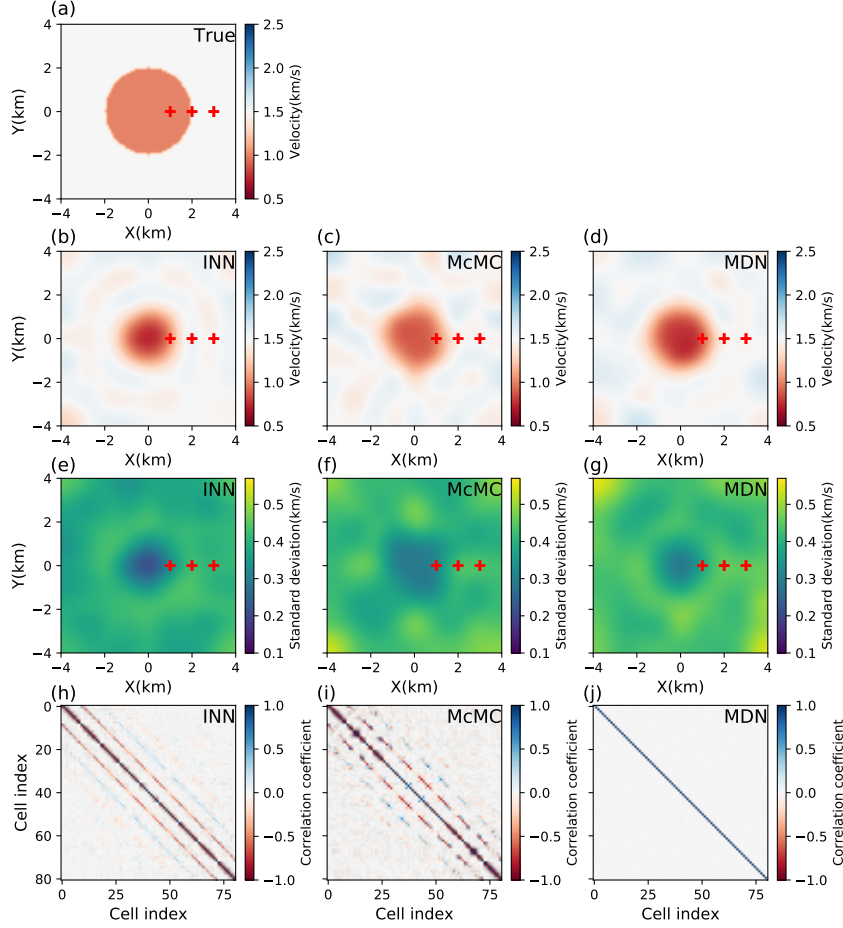


Figure 9. (a) The true velocity model. (b), (c) and (d) show the mean velocity models obtained using INN, McMC and MDN respectively. (e), (f) and (g) show the standard deviation at each point obtained using the three methods respectively. (h), (i) and (j) show the correlation coefficient matrices obtained using the three methods respectively. Red pluses show locations referred to in the text, at which marginal probability distributions are shown in Figure 10.

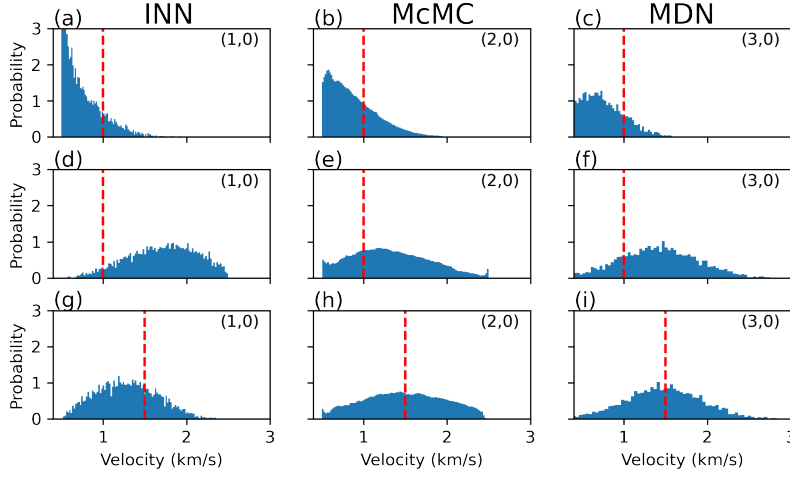


Figure 10. The marginal distributions at three points (red pluses in Figure 9) derived using different methods. (a), (b) and (c) show the marginal distributions at point (1,0) obtained using INN, McMC and MDN respectively. (d), (e) and (f) show similar marginal distributions at point (2,0), and (g), (h) and (i) show results for point (3,0). Red dashed line show the location of true value.

For the INN we use a network that contains eight reversible blocks (details in Appendix B2) trained using the ADAM optimizer (Kingma & Ba, 2014). To better understand the method, we compared the results with those obtained using MDNs and McMC. We train MDNs of the same network design as described in Earp and Curtis (2020) with 100 Gaussian kernels. Again the MDN network is trained 20 times with random initialization, and the network with best performance on test data is used to produce final results. For McMC we use a standard adaptive Metropolis-Hastings algorithm with a total of 6 chains, each of which contains 1,600,000 samples including a burn-in period of 600,000. Again the burn-in samples were ignored, and every 20th of the remaining 1,000,000 samples was included in the final set of McMC samples used for calculating statistics and marginal distributions.

As a first example we show results for data generated using a smooth velocity model which contains a low velocity anomaly in the centre within a homogeneous background (Figure 9a), similar to the test model in (Galetti et al., 2015; Zhang et al., 2018; Earp & Curtis, 2020; Zhang & Curtis, 2020a). This model is deliberately simple (which permits some degree of intuition about the tomographic solution), but the tomography problem is nevertheless substantially nonlinear (Galetti et al., 2015). The model also consists

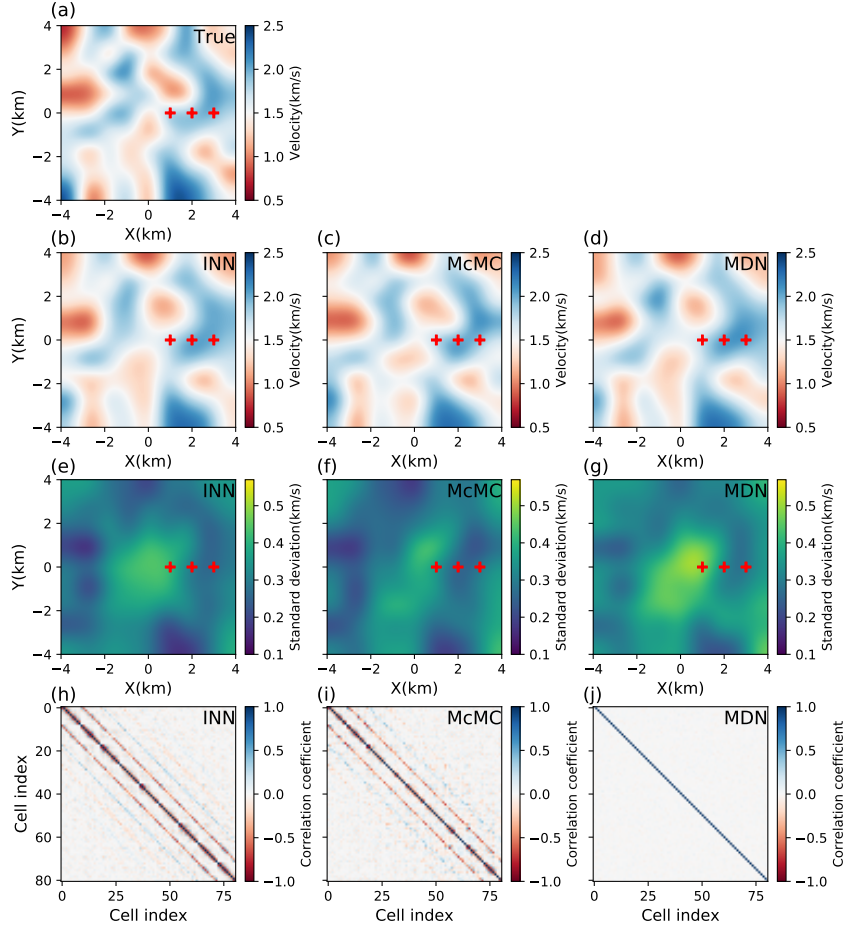


Figure 11. A random heterogeneous velocity example. Key as in Figure 9

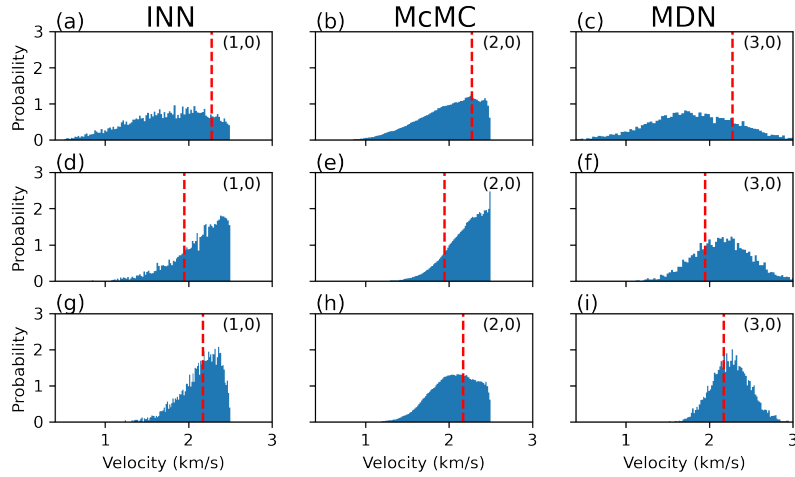


Figure 12. The marginal distributions at three points (red pluses in Figure 11) derived using different methods. Key as in Figure 10

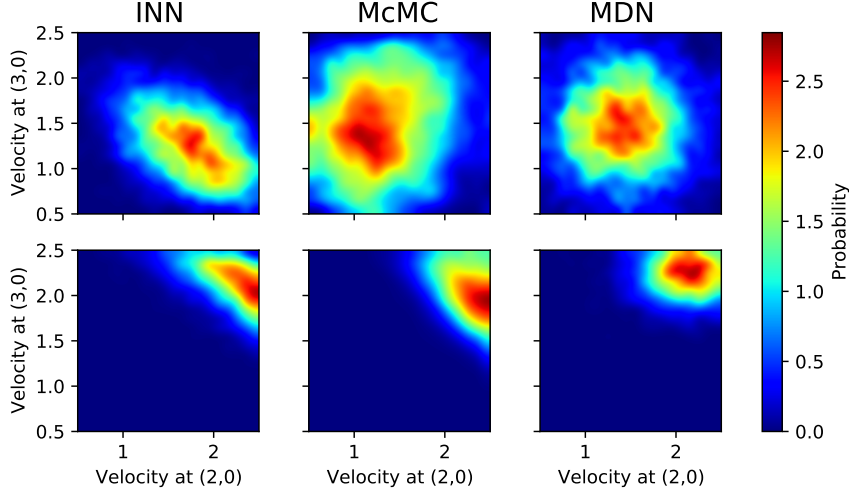


Figure 13. The bivariate marginal distributions between velocities at two locations: (2,0) km and (3,0) km obtained using INN, McMC and MDN for the two models in Figures 9 and 11 in upper and lower rows, respectively.

of both perfectly smooth regions and a sharp, spatially coherent boundary, whereas smooth regions are not represented by any randomly-selected training example (they are all of similar roughness to that in Figure 8a) so any reasonable solution must be interpolated between training examples. The example also presents the complication that the true model is defined on a relatively high resolution grid so that the circular boundary is not representable in the lower resolution parameterizations used for tomography (and for the set of training examples above). Zhang and Curtis (2020a) showed that the Bayesian solution to this problem varies significantly depending on the parameterization adopted, so in this example we use identical parameterizations for each of the three methods.

The travel time data are calculated using the high resolution grid (81×81), and are fed into the trained INNs and MDNs to predict posterior pdfs, and are also inverted using McMC to generate posterior samples. Overall the three methods produce similar mean (Figure 9b, c and d) and standard deviation models (Figure 9e, f and j). For example, all mean models show the middle low velocity anomaly, and small velocity variations around the anomaly which might be caused by lower resolution in those areas or by the random noise added to the data. Note also that we do not expect the mean model to match the true model. The mean is a statistic of the family of all models that might be true given the data; together with the standard deviation, it helps us to describe that

family (the goal of uncertainty analysis). Hence the mean is not an estimator of the true model, and particularly in nonlinear problems it is expected to deviate from the true model. All standard deviation models show low uncertainties across the central anomaly and higher uncertainties around the anomaly and in the four corners of the grid. Note that some detailed structure present in the results of McMC cannot be clearly observed in the results of INNs and MDNs. For example, there are four higher standard deviation anomalies above, below, left and right of the central anomaly, which are not clearly visible in the results of INNs and MDNs. This is probably because of cost function residuals remaining after training the INNs and MDNs.

Note that this is a high (81) dimensional problem, and there is minimal information in the uniformly-distributed prior pdf. The curse of dimensionality therefore implies that a huge number of samples would be needed to sample the parameter space adequately to explore all significant areas of the tomographic solution (Curtis & Lomax, 2001). Even sampling at only 2 samples per dimension would require $2^{81} = 2.4 \times 10^{24}$ samples. Therefore the training sets used are extremely small (far closer to 1 than to 2 samples per dimension), and even after 1.6 million samples the McMC method may not have converged to a statistically stable solution as it is difficult to assess convergence of McMC (Brooks et al., 2011). Therefore in this problem we do not know the exact Bayesian solution, nor which of the three solutions in Figure 9 is the most accurate. Nevertheless, the broad character of all three sets of means and standard deviations matches those found previously (papers cited above) so we have reasonable confidence in these statistics.

The correlation coefficient matrices obtained using INNs and McMC show similar results (Figure 9h and i). For example, there are negative correlations between neighbouring cells and positive correlations between every second neighbouring cells. In comparison the results obtained using MDNs do not show any correlation information. This is probably because MDNs become numerically unstable in high dimensionality and again use only a few kernels to represent the solution (Hjorth & Nabney, 1999; Rupprecht et al., 2017; Curro & Raquet, 2018; Cui et al., 2019; Makansi et al., 2019; Earp & Curtis, 2020).

To further analyse the results, Figure 10 shows marginal pdfs obtained by histogramming samples from the solutions of each of the three methods at three locations: (1,0), (2,0) and (3,0). Overall the results show similar marginal distributions, which suggests

that both INNs and MDNs can produce reasonable estimates of marginal distributions. However, the distributions obtained from different methods still show slightly different shapes which is probably caused by training residuals in INNs and MDNs. The upper row of Figure 13 shows the bivariate marginal distributions of velocities at location (2,0) and (3,0) obtained using the three methods. The results obtained using INN and McMC show clear correlation information between the two velocities; whereas those obtained using MDN do not. Similarly to the surface wave dispersion inversion examples, the distribution of normalized residuals obtained using INN and McMC are close to the standard Gaussian distribution (upper row in Figure S4); whereas that obtained using MDN is different. This again indicates that INNs can generate more accurate results than MDNs.

To explore generalization properties of the trained neural networks we show another example using data generated from a random velocity model in Figure 11a. Similarly to the first example, the travel time data are calculated using a denser grid (81×81). Overall the three methods produce similar mean and standard deviation models. The mean velocity models are largely the same as the true model. The standard deviation models show low uncertainties at the location of the high velocity anomaly at the east side of the area and at the location of the low velocity anomaly at the west side between $Y=0$ km and $Y=2$ km, and high uncertainties in the centre. Due to training residuals and possible lack of convergence of McMC the standard deviation models show other details that differ between methods. For example, the central high uncertainly anomaly adopts different shapes in the three results. Nevertheless, overall the results are fairly consistent.

The correlation coefficient matrices show similar results to the first example: the results from INNs and McMC show negative correlations between neighbouring cells and positive correlations between every second neighbouring cell, whereas the results from MDNs shows no correlation information. In Figure 12 we show marginal distributions at the same locations as for the first example. The results from MDNs are clearly approximately Gaussian, whereas the results from INNs show non-Gaussian shapes which provide more accurate approximations to results of McMC. Similarly to the first example, the bivariate marginal distributions (the bottom row in Figure 13) obtained using INN and McMC show clear trade-offs between the two parameters at the two locations; whereas those obtained using MDN failed to produce this information. Similarly to above the distributions of normalized residuals obtained using INN and McMC are close to the

standard Gaussian distribution (bottom row in Figure S4); whereas that obtained using MDN is different.

To demonstrate the effectiveness of the new loss function in equation 13, we train another INN using the same network design and same training procedure but with the loss function in equation 8. The results predicted by this INN are shown in the online supporting information (Figure S3 and S4). Although the mean models are very similar to those presented here, the standard deviation structure and the correlation matrices are significantly different. They failed to produce reliable estimates for uncertainties and correlation information between distant parameters. This clearly demonstrates the effectiveness of the new loss function.

Again the trained INNs and MDNs provide very efficient estimates of the posterior pdfs. For example, in a typical desktop the above trained MDN and INN takes about 3 seconds to produce a prediction of posterior pdfs – slightly longer than that required for 1D surface wave dispersion inversion because of the larger networks used here. In comparison McMC takes about 3 days on the same machine to generate the above results.

4 Application to field data

We apply the above method to ambient noise data of South England, which has complex geological structures compared to the rest of the British Isles (Nicolson et al., 2012, 2014; Galetti et al., 2017). We use ambient noise data that were recorded using 15 stations in 2006 to 2007 and in 2010 (Figure 14). The two horizontal components of the data (N and E) were first rotated to the transverse and radial directions, and the obtained transverse data were cross correlated to produce love waves between different station pairs. Travel times of group velocity at different periods between different station pairs were then estimated from those love waves. The details of the data processing procedure can be found in Galetti et al. (2017). In this example we study the group velocity map at 10 s.

For tomography we use a 9×10 regular grid of cells with a spacing of 0.4° in both latitude and longitude directions to parameterize the velocity model. The prior pdf for velocity is set to be a Uniform distribution between 1.8 km/s and 3.8 km/s, of which the lower and upper bound were chosen to be 0.5 km/s lower and higher than the minimum and maximum velocities between all available station pairs when assuming a great cir-

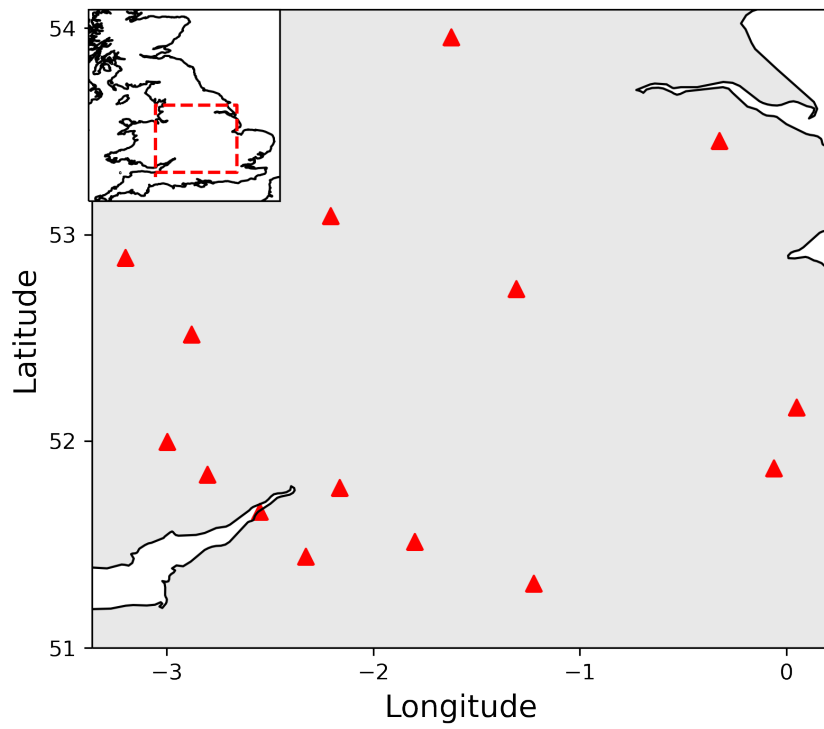


Figure 14. The receivers (red triangles) used in this study. The red dashed line box in the inset map denotes the location of the study area.

cle ray path. For likelihood function we use a Gaussian distribution, for which the data noise is estimated from independent travel time measurements by randomly stacking different subsets of daily cross correlations (Galetti et al., 2017).

For comparison we solve the tomographic problem using both INN and MCMC. For INN we use the same network design as used in the previous tomographic example and train the network in the same way. For training data we generate 200,000 velocity models from the prior pdf, and for each model the inter-station travel time data are calculated using the fast marching method. The corresponding noise for those travel time data is generated from Gaussian distributions whose standard deviation is set to be the value estimated from the data. Again 90 percent of them are used as training data and the remaining 10 percent are used as test data. The trained INN is then used to estimate the posterior pdf for the observed data using equation 4. For MCMC similarly to the previous example we run 6 independent chains, each of which contains 1,600,000 samples including a burn-in period of 600,000. The burn-in samples were discarded and every 20th of the remaining samples was collected as the final set of MCMC samples.

Figure 15 shows the mean, standard deviation and correlation coefficient matrix obtained using INN and MCMC respectively. Overall the two methods produce very similar results. For example, both results show a low velocity anomaly at latitude 53.0° and longitude -2.1° , which extends down to the south and connects with another low velocity anomaly. This anomaly has also been observed in the results obtained using reversible-jump MCMC (Galetti et al., 2017) and variational inference (Zhao et al., 2020), which is related to the Cheshire Basin. Another low velocity anomaly is also found in the south-east corresponding to the Anglian Basin. Both standard deviation models show lower uncertainties at the location of the first low velocity anomaly, and higher uncertainties at the location of the southeastern low velocity anomaly which is probably due to the lower data coverage across the latter area. Similarly there is a high velocity anomaly in the northeast associated with lower uncertainties. Compared to the synthetic example, the correlation coefficient matrices only show strong correlation between neighbouring cells (the red off-diagonal elements in Figure 15) which is probably caused by uneven data coverage and higher noise in the data.

In Figure 16 we show marginal distributions obtained using INN and MCMC at latitude 52.6 and longitude -3.0 , -2.2 , -1.4 and -0.6 (red pluses in Figure 15). For most of

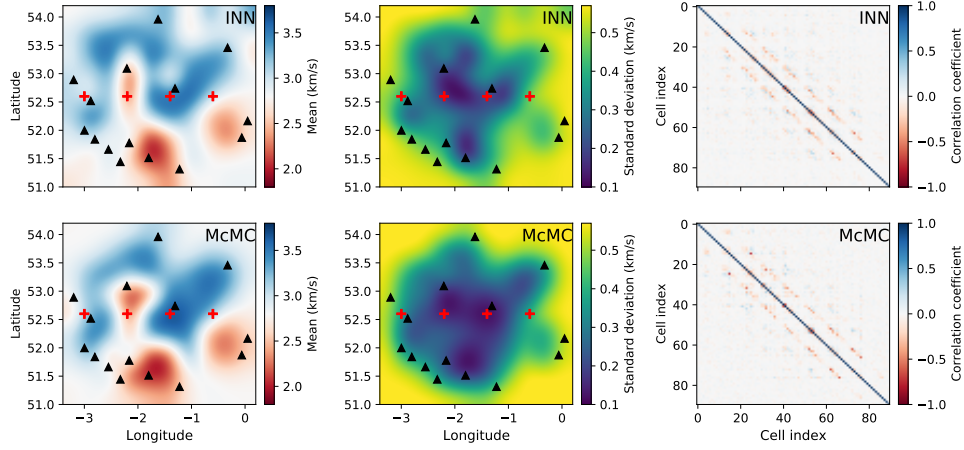


Figure 15. The mean, standard deviation and correlation matrix of the tomographic solution across southern England obtained using INN (**upper row**) and McMC (**bottom row**) respectively. The red pluses show the locations referred to in the text, at which marginal distributions are shown in Figure 16.

the locations the marginal distributions obtained using INN and McMC are very similar. At the location of longitude -1.4 the marginal distribution obtained using INN shows a slightly wider distribution compared to that obtained using McMC, which is probably caused by the training residual. This also explains the subtle difference presented in the standard deviation models obtained using INN and McMC.

We thus draw similar conclusion as in the synthetic example: the INN can provide an accurate approximation to the results obtained using McMC. Once the network is trained, it can produce predictions of the posterior pdf for new data very quickly. For example, the above example takes 3 seconds to produce a prediction on our desktop (32 Intel Xeon CPUs), whereas McMC takes about 40 hours to run on the same machine.

5 Discussion

Although trained neural networks provide efficient estimates of posterior pdfs, the methods require large numbers of training datasets to be created in advance, and training itself can still be computationally expensive. For example, on our desktop (32 Intel Xeon CPUs) it takes about 0.3 hours to generate 100,000 surface wave dispersion curves

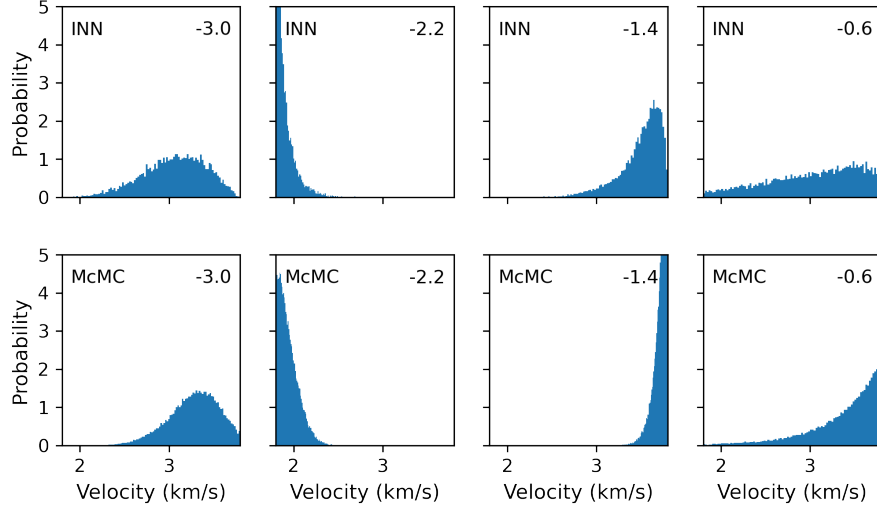


Figure 16. The marginal distributions at latitude 52.6 and longitude -3.0, -2.2, -1.4 and -0.6 (red pluses in Figure 15) obtained using INN (**upper row**) and McMC (**bottom row**) respectively.

Table 1. Comparison of computational cost required by different methods

Experiment	Surface wave inversion		Travel time tomography		Field data	
	CPU hours	GPU hours	CPU hours	GPU hours	CPU hours	GPU hours
INNs	0.3	6.0	3.3	19.0	2.1	13.5
MDN	0.3	0.25	3.3	3.0	NA	NA
McMC	4.5	NA	210.0	NA	120.0	NA

using one CPU core, 1.1 hours to generate 200,000 travel time data for the synthetic example and 0.7 hours for the field data example using 3 CPU cores. However the training data only need to be calculated once; even if additional prior information becomes available we can update the prior using the prior replacement method (Walker & Curtis, 2014) or the resampling method (Sambridge, 1999) rather than generating entirely new training samples. However, in that case although the prior samples can be augmented easily, retraining of the neural network on the new data is still required.

Neural networks can take between hours and days to train. For example, the above MDNs take 15 minutes to train for 1D surface wave dispersion inversion and 3 hours to train for the synthetic 2D travel time tomography example using one NVIDIA Tesla K80 GPU. In comparison, the above INNs take 6 hours to train for 1D surface wave disper-

sion inversion, and 19 hours to train for the synthetic travel time tomography example and 13.5 hours for the real data example using the same GPU. Table 1 summaries the computational cost of different methods for different tasks, including the training time (GPU hours) and data generation time (CPU hours). So although INNs improve the prediction accuracy of the posterior pdfs compared to MDN's, they take longer to train. This is because INNs are trained in both directions which generally requires larger networks to represent the forward and inverse process at the same time, and bidirectional training also intrinsically requires more computation and hence training time. However, training only needs to be done once: after training both types of neural networks can predict posterior pdfs in seconds. Computational efficiency is therefore gained when trained neural networks are applied many times, e.g., in real-time monitoring scenarios (Cao et al., 2020) or in highly parallelised (task-farmed) inference problems. Note that here we only compare the computational cost for one single training of INNs and MDNs, and for one single run of McMC. In practice both INNs and MDNs require additional training effort to tune the network architectures and training hyperparameters, and McMC also requires hyperparameter tuning, which makes comparison of the real time or cost required by each method difficult. Nevertheless the comparison presented here still provides valuable information about the different methods.

Even if the networks will only be applied once, it may be that they provide a more efficient solution than standard McMC, because of the generalisation property of neural networks. Essentially training networks is equivalent to performing a regression of the functional form of each network to the training data. Once that regression has been accomplished, the functional form approximately interpolates between examples in the training set, providing estimates of forward evaluations of infinitely many other model samples. This property implies that in some problems it may not be necessary to evaluate as many samples to train a neural network as would be required to characterise the solution using Monte Carlo or other purely sampling-based methods. However, note that this functional form is only valid within the prior range since there is no training data outside of the range of the prior distribution.

For example, in the above synthetic travel time tomography problems, the 1.6 million McMC samples required 3 days to evaluate, while the total times required to calculate training examples and train the MDN and INN were 4.3 hours and 20.1 hours, respectively. Hence in this example both the MDN and the INN were more efficient than

McMC even for one single inversion. Of course, the computational time strongly depends on hardware used in each case. In the above study different hardware is used for neural network training and McMC because McMC is difficult to parallelize, and hence it is difficult to take advantage of GPUs. It may also be the case that if we had used more efficient McMC methods such as Hamiltonian Monte Carlo (Duane et al., 1987; Neal et al., 2011; Fichtner et al., 2018) or Langevin Monte Carlo (Grenander & Miller, 1994) we could have improved the McMC performance to outstrip the two neural network methods. Nevertheless, this work demonstrates at least that neural network inversion can be competitive with Monte Carlo methods even for a single inversion, but additionally that they allow all subsequent inversions using the same prior information to be carried out almost for free.

Since INNs are trained bidirectionally, they also provide approximate forward functions. For example, for a given model \mathbf{m} and different data noise ϵ , the trained network can produce a distribution of data by feeding them into the network in the forward direction – this provides an approximation to the process described in equation 9. In Appendix C we show data distributions predicted by trained INNs for both surface wave dispersion curves and travel times for models described in section 3, and compared with solutions obtained by standard numerical modelling methods. The results show that trained INNs can also provide good approximations to the standard modelling methods. In this case since the standard forward modelling methods calculate dispersion curves and travel time data in seconds, INNs do not provide any benefits. However, for problems whose forward modelling is computationally expensive (e.g. 3D applications or full waveform modelling), INNs might provide a faster approximate forward function, and we note that neural networks have already been used to solve seismic forward modelling problems (Richardson, 2018; Song et al., 2020).

In this study we used independent Uniform distributions for each cell which can become ineffective for large inverse problems (Curtis & Lomax, 2001; de Pasquale & Linde, 2017; Earp & Curtis, 2020). If appropriate, the space of models that remain possible in the Bayesian solution can be reduced by using a smooth prior (Earp & Curtis, 2020) or structure-based priors (de Pasquale & Linde, 2017), or other more advanced priors can be used to improve computationally efficiency (Walker & Curtis, 2014; Zunino et al., 2015; de Pasquale & Linde, 2017; Ray & Myer, 2019; Caers, 2018; Mosser et al., 2020; M. A. Nawaz et al., 2020). We also used a fixed regular grid of cells to parametrize the subsurface, but

to increase flexibility we could use other parametrizations, such as Delauney triangulation (Curtis & Snieder, 1997), Voronoi diagrams (Sambridge et al., 1995; Bodin & Sambridge, 2009), wavelet representations (Fang et al., 2015; Zhang & Zhang, 2015) or other advanced parametrizations (Hawkins et al., 2019). The parametrization itself might also be predicted by neural networks along with the parameter values as in reversible-jump MCMC (Green, 1995; Bodin & Sambridge, 2009), for example, Gaussian processes, which can be treated as infinite, single layer neural networks (Neal, 2012), have been used with the reversible jump algorithm to predict both the parameterization and parameter values (Ray & Myer, 2019).

In this study we applied INNs to surface wave dispersion inversion and 2D travel time tomography, which usually do not show strong multimodality when using a fixed parameterization (Zhang & Curtis, 2020a; Zhao et al., 2020). However we note that theoretically INNs can be applied to predict any posterior distributions, including multimodal distributions. Therefore, future research would be to apply INNs to problems that have multimodal posterior distributions, for example, resistivity inversion (Ray et al., 2013) and full-waveform inversion (Zhang & Curtis, 2021).

INNs provide full posterior pdfs which might be impossible for very high dimensional problems, for example, full waveform electromagnetic or seismic inversion, because this may require a network that is too large to fit in memory and may need significantly more time to train. In this case one can train INNs to predict marginal distributions of only a few parameters, as was performed in Earp and Curtis (2020) in cases where dominant correlations were expected between parameters in some neighbourhood of each other. Advanced prior information which uses neural networks to reparameterize geologically correlated Earth structure with a lower number of parameters may also be used to reduce the dimensionality (Laloy et al., 2017; Mosser et al., 2020).

In this study we used coupling layers to implement invertible neural networks which might affect the expressiveness of the network. Other designs of invertible networks may be used as alternatives, for example, invertible residual networks (Behrmann et al., 2019) or Hamiltonian neural networks (Greydanus et al., 2019). Future research that makes a fair comparison between different architectures would be a useful contribution.

6 Conclusion

In this study we introduced invertible neural networks (INNs) to solve geophysical Bayesian inference problems. INNs are a class of neural networks that provide bijective mappings between inputs and outputs and can be trained to produce estimates of posterior probability density functions efficiently by introducing additional latent variables on the output (data) side. We applied the method to two types of problems: 1D surface wave dispersion inversion and 2D travel time tomography, and compared the results with those obtained using Markov chain Monte Carlo (McMC) and Mixture density networks (MDNs). The results show that INNs can provide accurate approximations of posterior pdfs obtained by McMC, including correlation information between parameters which is difficult to obtain using standard MDNs. The marginal distributions from INNs can also provide clearly non-Gaussian forms which are more similar to those obtained by McMC compared to the results obtained by MDNs. After training INNs can predict posterior pdfs in seconds, and therefore can be used to provide accurate estimates of posterior pdfs in rapid, real-time monitoring scenarios. Even accounting for training time, neural networks can be more efficient than Monte Carlo methods in applications to single inverse problems. It remains to be seen how far this latter result can be generalised to problems other than those tested here.

Acknowledgments

The authors thank the Edinburgh Imaging Project sponsors (BP, Schlumberger and Total) for supporting this research. This work has made use of the resources provided by the Edinburgh Compute and Data Facility (ECDF) (<http://www.ecdf.ed.ac.uk/>). The data and codes used in this study are available at Edinburgh DataShare (<https://datashare.ed.ac.uk/handle/10283/3802>).

References

- Adler, T. J., Ardizzone, L., Vemuri, A., Ayala, L., Gröhl, J., Kirchner, T., ... others (2019). Uncertainty-aware performance assessment of optical imaging modalities with invertible neural networks. *International journal of computer assisted radiology and surgery*, 14(6), 997–1007.
- Aki, K., & Lee, W. (1976). Determination of three-dimensional velocity anomalies under a seismic array using first P arrival times from local earthquakes: 1. a

- homogeneous initial model. *Journal of Geophysical research*, 81(23), 4381–4399.
- Araya-Polo, M., Jennings, J., Adler, A., & Dahlke, T. (2018). Deep-learning tomography. *The Leading Edge*, 37(1), 58–66.
- Ardizzone, L., Kruse, J., Wirkert, S., Rahner, D., Pellegrini, E. W., Klessen, R. S., ... Köthe, U. (2018). Analyzing inverse problems with invertible neural networks. *arXiv preprint arXiv:1808.04730*.
- Behrmann, J., Grathwohl, W., Chen, R. T., Duvenaud, D., & Jacobsen, J.-H. (2019). Invertible residual networks. In *International conference on machine learning* (pp. 573–582).
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. springer.
- Blei, D. M., Kucukelbir, A., & McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518), 859–877.
- Bodin, T., & Sambridge, M. (2009). Seismic tomography with the reversible jump algorithm. *Geophysical Journal International*, 178(3), 1411–1436.
- Bodin, T., Sambridge, M., Tkalčić, H., Arroucau, P., Gallagher, K., & Rawlinson, N. (2012). Transdimensional inversion of receiver functions and surface wave dispersion. *Journal of Geophysical Research: Solid Earth*, 117(B2).
- Bosch, M., Meza, R., Jiménez, R., & Hönig, A. (2006). Joint gravity and magnetic inversion in 3d using monte carlo methods. *Geophysics*, 71(4), G153–G156.
- Brocher, T. M. (2005). Empirical relations between elastic wavespeeds and density in the Earth’s crust. *Bulletin of the seismological Society of America*, 95(6), 2081–2092.
- Brooks, S., Gelman, A., Jones, G., & Meng, X.-L. (2011). *Handbook of Markov chain Monte Carlo*. CRC press.
- Caers, J. (2018). Bayesianism in the geosciences. In *Handbook of mathematical geosciences* (pp. 527–566). Springer, Cham.
- Cao, R., Earp, S., de Ridder, S. A., Curtis, A., & Galetti, E. (2020). Near-real-time near-surface 3d seismic velocity and uncertainty models by wavefield gradiometry and neural network inversion of ambient seismic noise. *Geophysics*, 85(1), KS13–KS27.
- Castagna, J. P., Batzle, M. L., & Eastwood, R. L. (1985). Relationships between

- compressional-wave and shear-wave velocities in clastic silicate rocks. *Geophysics*, 50(4), 571–581.
- Cui, H., Radosavljevic, V., Chou, F.-C., Lin, T.-H., Nguyen, T., Huang, T.-K., ... Djuric, N. (2019). Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In *2019 international conference on robotics and automation (icra)* (pp. 2090–2096).
- Curro, J., & Raquet, J. (2018). Deriving confidence from artificial neural networks for navigation. In *2018 IEEE/ION Position, Location and Navigation Symposium (Plans)* (pp. 1351–1361).
- Curtis, A., & Lomax, A. (2001). Prior information, sampling distributions, and the curse of dimensionality. *Geophysics*, 66(2), 372–378.
- Curtis, A., & Snieder, R. (1997). Reconditioning inverse problems using the genetic algorithm and revised parameterization. *Geophysics*, 62(5), 1524–1532.
- de Pasquale, G., & Linde, N. (2017). On structure-based priors in bayesian geophysical inversion. *Geophysical Journal International*, 208(3), 1342–1358.
- Devilee, R., Curtis, A., & Roy-Chowdhury, K. (1999). An efficient, probabilistic neural network approach to solving inverse problems: Inverting surface wave velocities for Eurasian crustal thickness. *Journal of Geophysical Research: Solid Earth*, 104(B12), 28841–28857.
- Devore, J. L., & Berk, K. N. (2012). *Modern mathematical statistics with applications*. Springer.
- de Wit, R. W., Valentine, A. P., & Trampert, J. (2013). Bayesian inference of earth’s radial seismic structure from body-wave traveltimes using neural networks. *Geophysical Journal International*, 195(1), 408–422.
- Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2016). Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*.
- Duane, S., Kennedy, A. D., Pendleton, B. J., & Roweth, D. (1987). Hybrid monte carlo. *Physics letters B*, 195(2), 216–222.
- Duputel, Z., Ferrazzini, V., Brenguier, F., Shapiro, N., Campillo, M., & Nercessian, A. (2009). Real time monitoring of relative velocity changes using ambient seismic noise at the piton de la fournaise volcano (la réunion) from january 2006 to june 2007. *Journal of Volcanology and Geothermal Research*, 184(1-2), 164–173.

- 836 Dzierwonski, A. M., & Woodhouse, J. H. (1987). Global images of the Earth's inte-
837 rior. *Science*, 236(4797), 37–48.
- 838 Earp, S., & Curtis, A. (2020). Probabilistic neural network-based 2d travel-time to-
839 mography. *NEURAL COMPUTING & APPLICATIONS*.
- 840 Earp, S., Curtis, A., Zhang, X., & Hansteen, F. (2020, 07). Probabilistic Neural
841 Network Tomography across Grane field (North Sea) from Surface Wave Dis-
842 persion Data. *Geophysical Journal International*. Retrieved from [https://](https://doi.org/10.1093/gji/ggaa328)
843 doi.org/10.1093/gji/ggaa328 (ggaa328) doi: 10.1093/gji/ggaa328
- 844 Fang, H., Yao, H., Zhang, H., Huang, Y.-C., & van der Hilst, R. D. (2015). Direct
845 inversion of surface wave dispersion for three-dimensional shallow crustal struc-
846 ture based on ray tracing: methodology and application. *Geophysical Journal*
847 *International*, 201(3), 1251–1263.
- 848 Fichtner, A., Zunino, A., & Gebraad, L. (2018). Hamiltonian monte carlo solution
849 of tomographic inverse problems. *Geophysical Journal International*, 216(2),
850 1344–1363.
- 851 Galetti, E., Curtis, A., Baptie, B., Jenkins, D., & Nicolson, H. (2017). Transdimen-
852 sional love-wave tomography of the British Isles and shear-velocity structure
853 of the east Irish Sea Basin from ambient-noise interferometry. *Geophysical*
854 *Journal International*, 208(1), 36–58.
- 855 Galetti, E., Curtis, A., Meles, G. A., & Baptie, B. (2015). Uncertainty loops in
856 travel-time tomography from nonlinear wave physics. *Physical review letters*,
857 114(14), 148501.
- 858 Gebraad, L., Boehm, C., & Fichtner, A. (2020). Bayesian elastic full-waveform
859 inversion using Hamiltonian Monte Carlo. *Journal of Geophysical Re-*
860 *search: Solid Earth*, 125(3), e2019JB018428. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019JB018428)
861 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019JB018428 doi:
862 10.1029/2019JB018428
- 863 Green, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and
864 Byesian model determination. *Biometrika*, 711–732.
- 865 Grenander, U., & Miller, M. I. (1994). Representations of knowledge in complex
866 systems. *Journal of the Royal Statistical Society: Series B (Methodological)*,
867 56(4), 549–581.
- 868 Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., & Smola, A. (2012).

- 869 A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1),
870 723–773.
- 871 Greydanus, S., Dzamba, M., & Yosinski, J. (2019). Hamiltonian neural networks. In
872 *Advances in neural information processing systems* (pp. 15379–15389).
- 873 Haario, H., Saksman, E., Tamminen, J., et al. (2001). An adaptive metropolis algo-
874 rithm. *Bernoulli*, 7(2), 223–242.
- 875 Hawkins, R., Bodin, T., Sambridge, M., Choblet, G., & Husson, L. (2019). Trans-
876 dimensional surface reconstruction with different classes of parameterization.
877 *Geochemistry, Geophysics, Geosystems*, 20(1), 505–529.
- 878 Herrmann, R. B. (2013). Computer programs in seismology: An evolving tool for in-
879 struction and research. *Seismological Research Letters*, 84(6), 1081–1088.
- 880 Hjorth, L. U., & Nabney, I. T. (1999). Regularisation of mixture density networks.
881 In *1999 ninth international conference on artificial neural networks icann*
882 *99.(conf. publ. no. 470)* (Vol. 2, pp. 521–526).
- 883 Huang, G., Nammour, R., & Symes, W. W. (2018). Volume source-based extended
884 waveform inversion. *Geophysics*, 83(5), R369–R387.
- 885 Iyer, H., & Hirahara, K. (1993). *Seismic tomography: Theory and practice*. Springer
886 Science & Business Media.
- 887 Jan van Heijst, H., Snieder, R., & Nowack, R. (1994). Resolving a low-velocity zone
888 with surface-wave data. *Geophysical Journal International*, 118(2), 333–343.
- 889 Jones, A. G., & Hutton, R. (1979). A multi-station magnetotelluric study in south-
890 ern scotland—ii. monte-carlo inversion of the data and its geophysical and
891 tectonic implications. *Geophysical Journal International*, 56(2), 351–368.
- 892 Karmakar, M., & Maiti, S. (2019). Short term memory efficient pore pressure pre-
893 diction via bayesian neural networks at bering sea slope of iodp expedition 323.
894 *Measurement*, 135, 852–868.
- 895 Karmakar, M., Maiti, S., Singh, A., Ojha, M., & Maity, B. S. (2018). Mapping
896 of rock types using a joint approach by combining the multivariate statistics,
897 self-organizing map and bayesian neural networks: an example from iodp 323
898 site. *Marine Geophysical Research*, 39(3), 407–419.
- 899 Käufel, P., P. Valentine, A., W. de Wit, R., & Trampert, J. (2016). Solving proba-
900 bilistic inverse problems rapidly with prior samples. *Geophysical Journal Inter-*
901 *national*, 205(3), 1710–1728.

- 902 Käüfl, P., Valentine, A., de Wit, R., & Trampert, J. (2015). Robust and fast prob-
 903 abilistic source parameter estimation from near-field displacement waveforms
 904 using pattern recognition. *Bulletin of the Seismological Society of America*,
 905 *105*(4), 2299–2312.
- 906 Käüfl, P., Valentine, A. P., O’Toole, T. B., & Trampert, J. (2014). A framework for
 907 fast probabilistic centroid-moment-tensor determination—inversion of regional
 908 static displacement measurements. *Geophysical Journal International*, *196*(3),
 909 1676–1693.
- 910 Khoshkholgh, S., Zunino, A., & Mosegaard, K. (2020). Informed proposal monte
 911 carlo. *arXiv preprint arXiv:2005.14398*.
- 912 Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv*
 913 *preprint arXiv:1412.6980*.
- 914 Kingma, D. P., & Dhariwal, P. (2018). Glow: Generative flow with invertible 1x1
 915 convolutions. In *Advances in neural information processing systems* (pp.
 916 10215–10224).
- 917 Koohababni, N. A., Jahanifar, M., Gooya, A., & Rajpoot, N. (2018). Nuclei de-
 918 tection using mixture density networks. In *International workshop on machine*
 919 *learning in medical imaging* (pp. 241–248).
- 920 Kotsi, M., Malcolm, A., & Ely, G. (2020). Time-lapse full-waveform inversion us-
 921 ing hamiltonian monte carlo: A proof of concept. In *Seg technical program ex-*
 922 *panded abstracts 2020* (pp. 845–849). Society of Exploration Geophysicists.
- 923 Laloy, E., Hérault, R., Lee, J., Jacques, D., & Linde, N. (2017). Inversion using a
 924 new low-dimensional representation of complex binary geological media based
 925 on a deep neural network. *Advances in water resources*, *110*, 387–405.
- 926 Laloy, E., Linde, N., & Jacques, D. (2019). Approaching geoscientific inverse prob-
 927 lems with adversarial vector-to-image domain transfer networks. *arXiv preprint*
 928 *arXiv:1912.09954*.
- 929 Makansi, O., Ilg, E., Cicek, O., & Brox, T. (2019). Overcoming limitations of
 930 mixture density networks: A sampling and fitting framework for multimodal
 931 future prediction. In *Proceedings of the ieee conference on computer vision and*
 932 *pattern recognition* (pp. 7144–7153).
- 933 Malinverno, A., Leaney, S., et al. (2000). A Monte Carlo method to quantify uncer-
 934 tainty in the inversion of zero-offset VSP data. In *2000 seg annual meeting*.

- 935 Meier, U., Curtis, A., & Trampert, J. (2007a). A global crustal model constrained
936 by nonlinearised inversion of fundamental mode surface waves. *Geophysical Re-*
937 *search Letters*, *34*, L16304.
- 938 Meier, U., Curtis, A., & Trampert, J. (2007b). Global crustal thickness from neu-
939 ral network inversion of surface wave data. *Geophysical Journal International*,
940 *169*(2), 706–722.
- 941 Men, Z., Yee, E., Lien, F.-S., Wen, D., & Chen, Y. (2016). Short-term wind speed
942 and power forecasting using an ensemble of mixture density neural networks.
943 *Renewable Energy*, *87*, 203–211.
- 944 Minsley, B. J. (2011). A trans-dimensional bayesian markov chain monte carlo algo-
945 rithm for model assessment using frequency-domain electromagnetic data. *Geo-*
946 *physical Journal International*, *187*(1), 252–272.
- 947 Mohamed, L., Christie, M., & Demyanov, V. (2010). Comparison of stochastic sam-
948 pling algorithms for uncertainty quantification. *SPE Journal*, *15*(01), 31–38.
- 949 Moran, O., Caramazza, P., Faccio, D., & Murray-Smith, R. (2018). Deep, complex,
950 invertible networks for inversion of transmission effects in multimode optical
951 fibres. In *Advances in neural information processing systems* (pp. 3280–3291).
- 952 Mosegaard, K., & Tarantola, A. (1995). Monte Carlo sampling of solutions to
953 inverse problems. *Journal of Geophysical Research: Solid Earth*, *100*(B7),
954 12431–12447.
- 955 Mosser, L., Dubrule, O., & Blunt, M. J. (2020). Stochastic seismic waveform inver-
956 sion using generative adversarial networks as a geological prior. *Mathematical*
957 *Geosciences*, *52*(1), 53–79.
- 958 Moya, A., & Irikura, K. (2010). Inversion of a velocity model using artificial neural
959 networks. *Computers & geosciences*, *36*(12), 1474–1483.
- 960 Nawaz, M., & Curtis, A. (2019). Rapid discriminative variational Bayesian inversion
961 of geophysical data for the spatial distribution of geological properties. *Journal*
962 *of Geophysical Research: Solid Earth*.
- 963 Nawaz, M. A., & Curtis, A. (2018). Variational Bayesian inversion (VBI) of quasi-
964 localized seismic attributes for the spatial distribution of geological facies. *Geo-*
965 *physical Journal International*, *214*(2), 845–875.
- 966 Nawaz, M. A., Curtis, A., Shahraeeni, M. S., & Gerea, C. (2020). Variational
967 Bayesian inversion of seismic attributes jointly for geological facies and

- petrophysical rock properties. *GEOPHYSICS*, 1-78. Retrieved from
<https://doi.org/10.1190/geo2019-0163.1> doi: 10.1190/geo2019-0163.1
- Neal, R. M. (2012). *Bayesian learning for neural networks* (Vol. 118). Springer Science & Business Media.
- Neal, R. M., et al. (2011). Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*, 2(11), 2.
- Nicolson, H., Curtis, A., & Baptie, B. (2014). Rayleigh wave tomography of the British Isles from ambient seismic noise. *Geophysical Journal International*, 198(2), 637–655.
- Nicolson, H., Curtis, A., Baptie, B., & Galetti, E. (2012). Seismic interferometry and ambient noise tomography in the British Isles. *Proceedings of the Geologists' Association*, 123(1), 74–86.
- Osborne, C. M., Armstrong, J. A., & Fletcher, L. (2019). Radynversion: learning to invert a solar flare atmosphere with invertible neural networks. *The Astrophysical Journal*, 873(2), 128.
- Piana Agostinetti, N., Giacomuzzi, G., & Malinverno, A. (2015). Local three-dimensional earthquake tomography by trans-dimensional Monte Carlo sampling. *Geophysical Journal International*, 201(3), 1598–1617.
- Ramdas, A., Reddi, S. J., Póczos, B., Singh, A., & Wasserman, L. (2015). On the decreasing power of kernel and distance based nonparametric hypothesis tests in high dimensions. In *Twenty-ninth aaai conference on artificial intelligence*.
- Rawlinson, N., & Sambridge, M. (2004). Multiple reflection and transmission phases in complex layered media using a multistage fast marching method. *Geophysics*, 69(5), 1338–1350.
- Ray, A., Alumbaugh, D. L., Hoversten, G. M., & Key, K. (2013). Robust and accelerated Bayesian inversion of marine controlled-source electromagnetic data using parallel tempering. *Geophysics*, 78(6), E271–E280.
- Ray, A., Kaplan, S., Washbourne, J., & Albertin, U. (2017). Low frequency full waveform seismic inversion within a tree based Bayesian framework. *Geophysical Journal International*, 212(1), 522–542.
- Ray, A., & Myer, D. (2019). Bayesian geophysical inversion with trans-dimensional gaussian process machine learning. *Geophysical Journal International*, 217(3), 1706–1726.

- 1001 Ray, A., Sekar, A., Hoversten, G. M., & Albertin, U. (2016). Frequency domain
1002 full waveform elastic inversion of marine seismic data from the Alba field using
1003 a Bayesian trans-dimensional algorithm. *Geophysical Journal International*,
1004 205(2), 915–937.
- 1005 Richardson, A. (2018). Seismic full-waveform inversion using deep learning tools and
1006 techniques. *arXiv preprint arXiv:1801.07232*.
- 1007 Richmond, K. (2007). Trajectory mixture density networks with multiple mix-
1008 tures for acoustic-articulatory inversion. In *International conference on nonlin-*
1009 *ear speech processing* (pp. 263–272).
- 1010 Rizzuti, G., Siahkoohi, A., Witte, P. A., & Herrmann, F. J. (2020). Parameterizing
1011 uncertainty by deep invertible networks: An application to reservoir charac-
1012 terization. In *Seg technical program expanded abstracts 2020* (pp. 1541–1545).
1013 Society of Exploration Geophysicists.
- 1014 Rossi, L. (2017). Bayesian gravity inversion by monte carlo methods.
- 1015 Röth, G., & Tarantola, A. (1994). Neural networks and inversion of seismic data.
1016 *Journal of Geophysical Research: Solid Earth*, 99(B4), 6753–6768.
- 1017 Rupprecht, C., Laina, I., DiPietro, R., Baust, M., Tombari, F., Navab, N., & Hager,
1018 G. D. (2017). Learning in an uncertain world: Representing ambiguity through
1019 multiple hypotheses. In *Proceedings of the ieee international conference on*
1020 *computer vision* (pp. 3591–3600).
- 1021 Sahin, G. G., & Gurevych, I. (2020). Two birds with one stone: Investigating invert-
1022 ible neural networks for inverse problems in morphology. In *Aaai* (pp. 7814–
1023 7821).
- 1024 Salvatier, J., Wiecki, T. V., & Fonnesbeck, C. (2016). Probabilistic programming in
1025 python using pymc3. *PeerJ Computer Science*, 2, e55.
- 1026 Sambridge, M. (1999). Geophysical inversion with a neighbourhood algorithm – i.
1027 searching a parameter space. *Geophysical journal international*, 138(2), 479–
1028 494.
- 1029 Sambridge, M., Braun, J., & McQueen, H. (1995). Geophysical parametrization and
1030 interpolation of irregular data using natural neighbours. *Geophysical Journal*
1031 *International*, 122(3), 837–857.
- 1032 Shahraeeni, M. S., & Curtis, A. (2011). Fast probabilistic nonlinear petrophysical in-
1033 version. *Geophysics*, 76(2), E45–E58.

- Shahraeeni, M. S., Curtis, A., & Chao, G. (2012). Fast probabilistic petrophysical mapping of reservoirs from 3d seismic data. *Geophysics*, 77(3), O1–O19.
- Shapiro, N. M., Campillo, M., Stehly, L., & Ritzwoller, M. H. (2005). High-resolution surface-wave tomography from ambient seismic noise. *Science*, 307(5715), 1615–1618.
- Shen, W., Ritzwoller, M. H., Schulte-Pelkum, V., & Lin, F.-C. (2012). Joint inversion of surface wave dispersion and receiver functions: a Bayesian Monte-Carlo approach. *Geophysical Journal International*, 192(2), 807–836.
- Siahkoobi, A., Rizzuti, G., Witte, P. A., & Herrmann, F. J. (2020). Faster uncertainty quantification for inverse problems with conditional normalizing flows. *arXiv preprint arXiv:2007.07985*.
- Smith, R. C. (2013). *Uncertainty quantification: theory, implementation, and applications* (Vol. 12). Siam.
- Song, C., Alkhalifah, T., & Waheed, U. b. (2020). Solving the acoustic vti wave equation using physics-informed neural networks. *arXiv preprint arXiv:2008.01865*.
- Subbey, S., Mike, C., Sambridge, M., et al. (2003). A strategy for rapid quantification of uncertainty in reservoir performance prediction. In *Spe reservoir simulation symposium*.
- Tolstikhin, I., Bousquet, O., Gelly, S., & Schoelkopf, B. (2017). Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*.
- Walker, M., & Curtis, A. (2014). Varying prior information in bayesian inversion. *Inverse Problems*, 30(6), 065002.
- Williams, P. M. (1996). Using neural networks to model conditional multivariate densities. *Neural Computation*, 8(4), 843–854.
- You, N., Li, Y. E., & Cheng, A. (2020). Shale anisotropy model building based on deep neural networks. *Journal of Geophysical Research: Solid Earth*, 125(2), e2019JB019042.
- Young, M. K., Rawlinson, N., & Bodin, T. (2013). Transdimensional inversion of ambient seismic noise for 3D shear velocity structure of the Tasmanian crust. *Geophysics*, 78(3), WB49–WB62.
- Zhang, X., & Curtis, A. (2020a). Seismic tomography using variational inference methods. *Journal of Geophysical Research: Solid Earth*, 125(4),

- 1067 e2019JB018589.
- 1068 Zhang, X., & Curtis, A. (2020b). Variational full-waveform inversion. *Geophysical*
 1069 *Journal International*, 222(1), 406–411.
- 1070 Zhang, X., & Curtis, A. (2021). Bayesian full-waveform inversion with realistic pri-
 1071 ors. *arXiv preprint arXiv:2104.04775*.
- 1072 Zhang, X., Curtis, A., Galetti, E., & de Ridder, S. (2018). 3-D Monte Carlo surface
 1073 wave tomography. *Geophysical Journal International*, 215(3), 1644–1658.
- 1074 Zhang, X., Hansteen, F., Curtis, A., & de Ridder, S. (2020). 1D, 2D and 3D Monte
 1075 Carlo ambient noise tomography using a dense passive seismic array installed
 1076 on the north sea seabed. *Journal of Geophysical Research: Solid Earth*, 125(2),
 1077 e2019JB018552. doi: 10.1029/2019JB018552
- 1078 Zhang, X., Roy, C., Curtis, A., Nowacki, A., & Baptie, B. (2020, 05). Imaging
 1079 the subsurface using induced seismicity and ambient noise: 3-D tomographic
 1080 Monte Carlo joint inversion of earthquake body wave traveltimes and surface
 1081 wave dispersion. *Geophysical Journal International*, 222(3), 1639–1655. doi:
 1082 10.1093/gji/ggaa230
- 1083 Zhang, X., & Zhang, H. (2015). Wavelet-based time-dependent travel time tomog-
 1084 raphy method and its application in imaging the Etna volcano in Italy. *Journal*
 1085 *of Geophysical Research: Solid Earth*, 120(10), 7068–7084.
- 1086 Zhao, X., Curtis, A., & Zhang, X. (2020). Bayesian seismic tomography using nor-
 1087 malizing flows. *EarthArXiv*.
- 1088 Zunino, A., Mosegaard, K., Lange, K., Melnikova, Y., & Mejer Hansen, T. (2015).
 1089 Monte carlo reservoir analysis combining seismic reflection data and informed
 1090 priors. *Geophysics*, 80(1), R31–R41.

1091 **Appendix A Proof of convergence**

1092 We follow the proof of Ardizzone et al. (2018) to prove the convergence of the method
 1093 described in section 2.3. Denote the prior probability distribution of data as $p(\mathbf{d})$, the
 1094 chosen probability distribution of latent variables as $p(\mathbf{z})$ and the joint probability dis-
 1095 tribution of network output as $q(\mathbf{d}, \mathbf{z})$. For a given measurement \mathbf{d}^* , define the poste-
 1096 rior distribution $p(\mathbf{m}|\mathbf{d}^*)$ as the distribution of models $\{\mathbf{m}\}$ which generate the exact
 1097 data \mathbf{d}^* .

Theorem: If an INN $f(\mathbf{m}) = [\mathbf{d}, \mathbf{z}]$ is trained such that both the loss function $\mathcal{L}_d = \|\mathbf{d} - f_d(\mathbf{m})\|$ and $\mathcal{L}_z = \text{MMD}[q(\mathbf{d}, \mathbf{z}), p(\mathbf{d})p(\mathbf{z})]$ equal zero, sampling according to equation 4 produces the posterior $p(\mathbf{m}|\mathbf{d}^*)$ for a given measurement \mathbf{d}^* .

Proof: Because MMD only equals zero when the two distributions are identical (Gretton et al., 2012), the loss function $\mathcal{L}_z = 0$ implies that:

$$q(\mathbf{d}, \mathbf{z}) = p(\mathbf{d})p(\mathbf{z}) \quad (\text{A1})$$

Suppose we take the posterior distribution $p(\mathbf{m}|\mathbf{d}^*)$ for a given \mathbf{d}^* , and transform it through the forward process of our perfectly trained INN. This results in a joint distribution $q^*(\mathbf{d}, \mathbf{z})$. According to the definition of $p(\mathbf{m}|\mathbf{d}^*)$ and the fact that $\mathcal{L}_d = 0$, the distribution of network output \mathbf{d} (marginalized over \mathbf{z}) is $q^*(\mathbf{d}) = \delta(\mathbf{d} - \mathbf{d}^*)$. Due to the fact that $\mathcal{L}_z = 0$, \mathbf{d} and \mathbf{z} are independent and the output distribution of \mathbf{z} is still $p(\mathbf{z})$. The distribution $q^*(\mathbf{d}, \mathbf{z})$ can therefore be expressed as

$$q^*(\mathbf{d}, \mathbf{z}) = \delta(\mathbf{d} - \mathbf{d}^*)p(\mathbf{z}) \quad (\text{A2})$$

This implies that we can transform the distribution $\delta(\mathbf{d} - \mathbf{d}^*)p(\mathbf{z})$ to the posterior $p(\mathbf{m}|\mathbf{d}^*)$ through the backward process of the INN. Equivalently speaking, we can repeatedly input \mathbf{d}^* with randomly sampled \mathbf{z} from $p(\mathbf{z})$ to the backward process of the INN, and obtain the posterior distribution $p(\mathbf{m}|\mathbf{d}^*)$.

Note that the above proof relies on the definition that the posterior distribution $p(\mathbf{m}|\mathbf{d}^*)$ is the distribution of models that generate the exact data \mathbf{d}^* . In Bayes theorem this means that the likelihood function is $\delta(\mathbf{d} - \mathbf{d}^*)$, which is not true in reality as observed data always contain noise. This issue can be resolved by treating noise as a part of the forward modelling process, that is,

$$\mathbf{d} = \mathbf{d}_{sim} + \epsilon = F(\mathbf{m}) + \epsilon \quad (\text{A3})$$

where \mathbf{d} represents the observed, noisy data, \mathbf{d}_{sim} is the simulated data using forward function $F(\mathbf{m})$, and ϵ is random noise generated from some distribution. In this way the method, and the above proof, can be invoked without any change.

Table B1. Training data and model and data dimensionalities used in different experiments

Experiment	training data	dim(inputs)	dim(d)	dim(z)
Surface wave inversion	100,000	24	14	22
Travel time tomography	200,000	241	120	181
Field data	200,000	121	31	99

Appendix B Network configuration

Table B1 summarizes the training datasets and model and data dimensionalities used in 1D surface wave dispersion inversion and 2D travel time tomography. Note that in practice zeros are padded to the inputs to ensure the same dimensionality.

B1 Network configuration for surface wave dispersion inversion

INN: 4 reversible blocks, each of which contains two coupling layers as described in equation 2. Each affine function (i.e. s_i and t_i) is implemented using a neural network with 3 fully connected layers each of which contains 512 hidden units with RELU activation functions. The ADAM optimizer is used with a batch size of 1000.

MDN: 20 mixture Gaussian kernels, 4 fully connected layers with RELU activation functions. The sizes of layers are 200, 300, 200 and 200 respectively. The ADAM optimizer is used with a batch size of 1000.

B2 Network configuration for travel time tomography

INN: 8 reversible blocks as described in equation 2. Each affine function (i.e. s_i and t_i) is implemented using one convolutional layer with 32 channels for the first four blocks, and one fully connected layer containing 1024 hidden units for the remaining blocks. The ADAM optimizer is used with a batch size of 1000.

MDN: 100 mixture Gaussian kernels, 7 layers in total containing 3 convolutional layers and 4 fully connected layers. The number of channels of the 3 convolutional layers are 128, 128 and 64 respectively. The size of the 4 fully connected layers are 800, 150, 600 and 1500 respectively. The ADAM optimizer is used with a batch size of 1000.

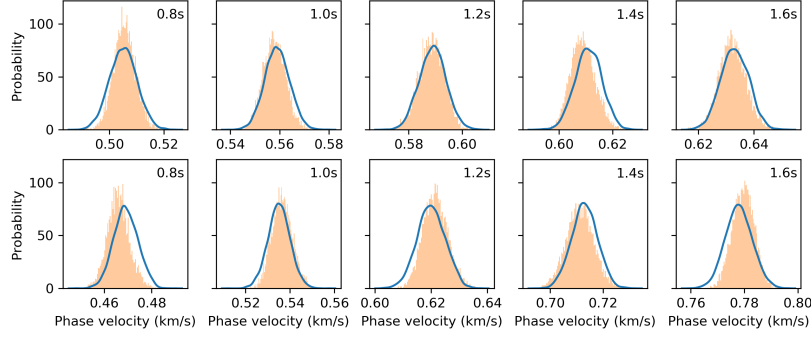


Figure C1. Distributions of phase velocities with random noise predicted by INNs (orange histograms) for a specific shear velocity model at periods 0.8 s, 1.0 s, 1.2 s, 1.4 s and 1.6 s. The distributions in upper and lower rows correspond to the upper and lower models respectively, shown by red lines in Figure 5. Blue lines show true distributions obtained using the standard forward modelling method with random noise added to the synthetic data.

Appendix C Approximate forward modelling function using INNs

Since INNs are trained bidirectionally, they also provide approximate forward functions. For example, one can obtain the distribution of data with noise for a fixed model \mathbf{m} by running the network forward with noise parameter ϵ (Figure 1) distributed according to its assumed distribution. Figure C1 shows phase velocity distributions obtained from INNs for those shear velocity models used in the surface wave dispersion inversion experiment at 5 different periods, and compares the results with those obtained using standard forward modelling methods using equation 9, which we refer to as true distributions. The noise of each phase velocity is assumed to follow a Gaussian distribution with a standard deviation of 5 m/s. The results show that the distributions predicted by INNs from the training samples alone, can provide good approximations to those obtained by the forward models themselves. Note that because of training residuals the distributions obtained from INNs are slightly different from the true distributions.

Similarly the upper and lower rows in Figure C2 show distributions of travel times for 5 randomly selected (virtual) source-receiver pairs predicted by INNs for the two velocity models in Figure 9 and Figure 11 respectively. The noise of each travel time is assumed to follow a Gaussian distribution with a standard deviation of 0.05 s. The results show that distributions obtained from INNs are largely similar to the true distributions

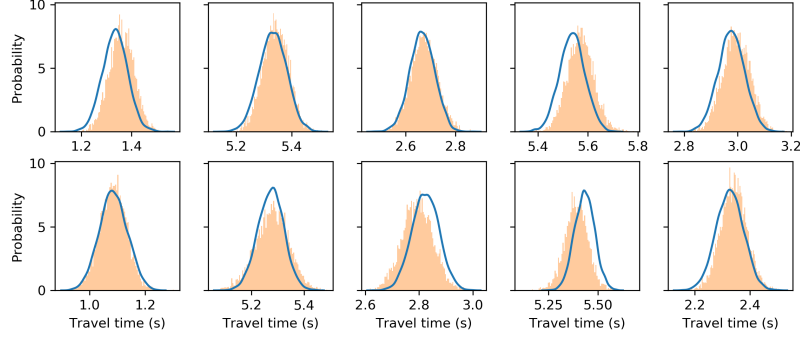


Figure C2. Travel time distributions for 5 randomly selected (virtual) source-receiver pairs predicted by INNs (orange histograms) for (**upper row**) the smooth velocity model in Figure 9a and (**lower row**) the random velocity model in Figure 11a. Blue lines show true distributions obtained using the standard forward modelling method with random noise added to the synthetic data.

1163 (those obtained using standard forward modelling method). However, there are still dif-
 1164 ferences between the two distributions caused by training residuals of INNs.